

Архитектура ВКонтакте: там, где данные

Илья Щербак



О себе



Родился



Учился



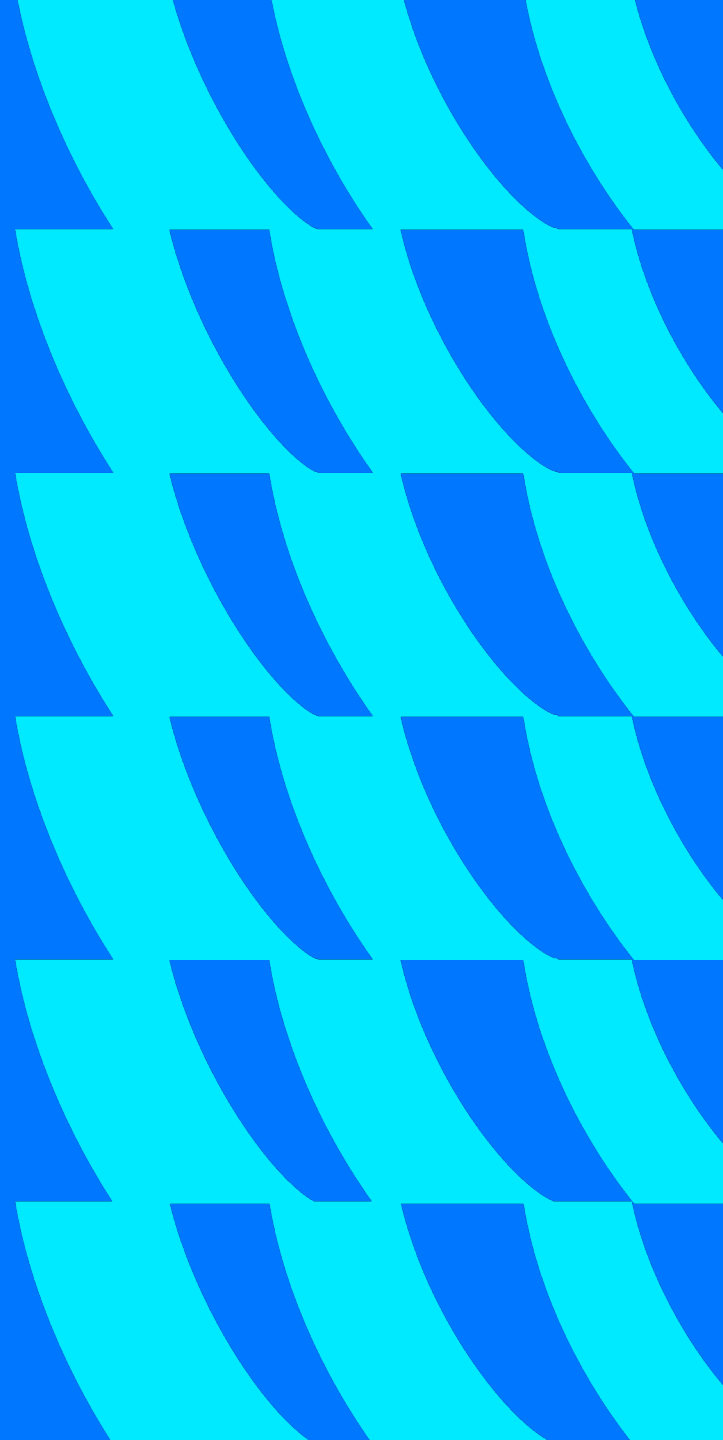
Чему-то
научился



Работаю
ВКонтакте



Команда Баз данных и инфраструктуры ВКонтакте







Технический вектор команды



Отказоустойчивость



Быстродействие



Скорость разработки
продукта

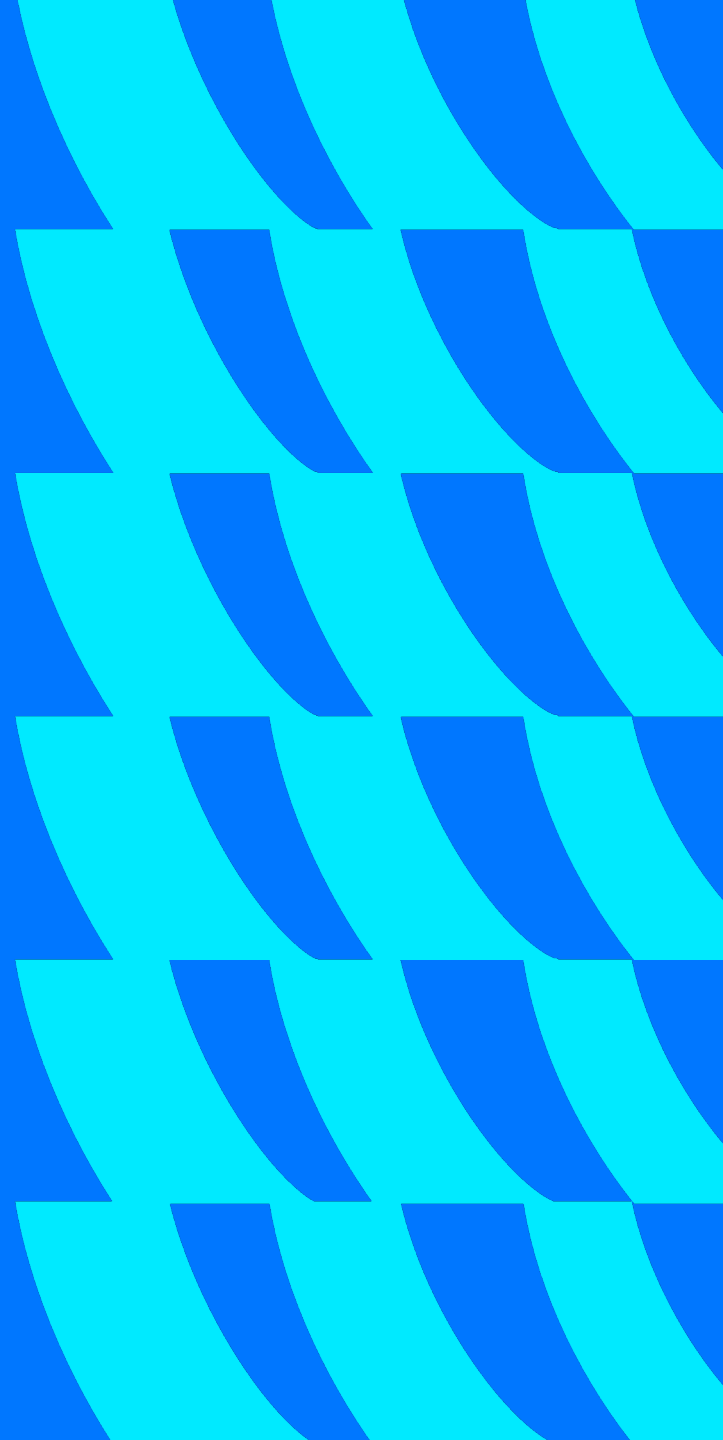


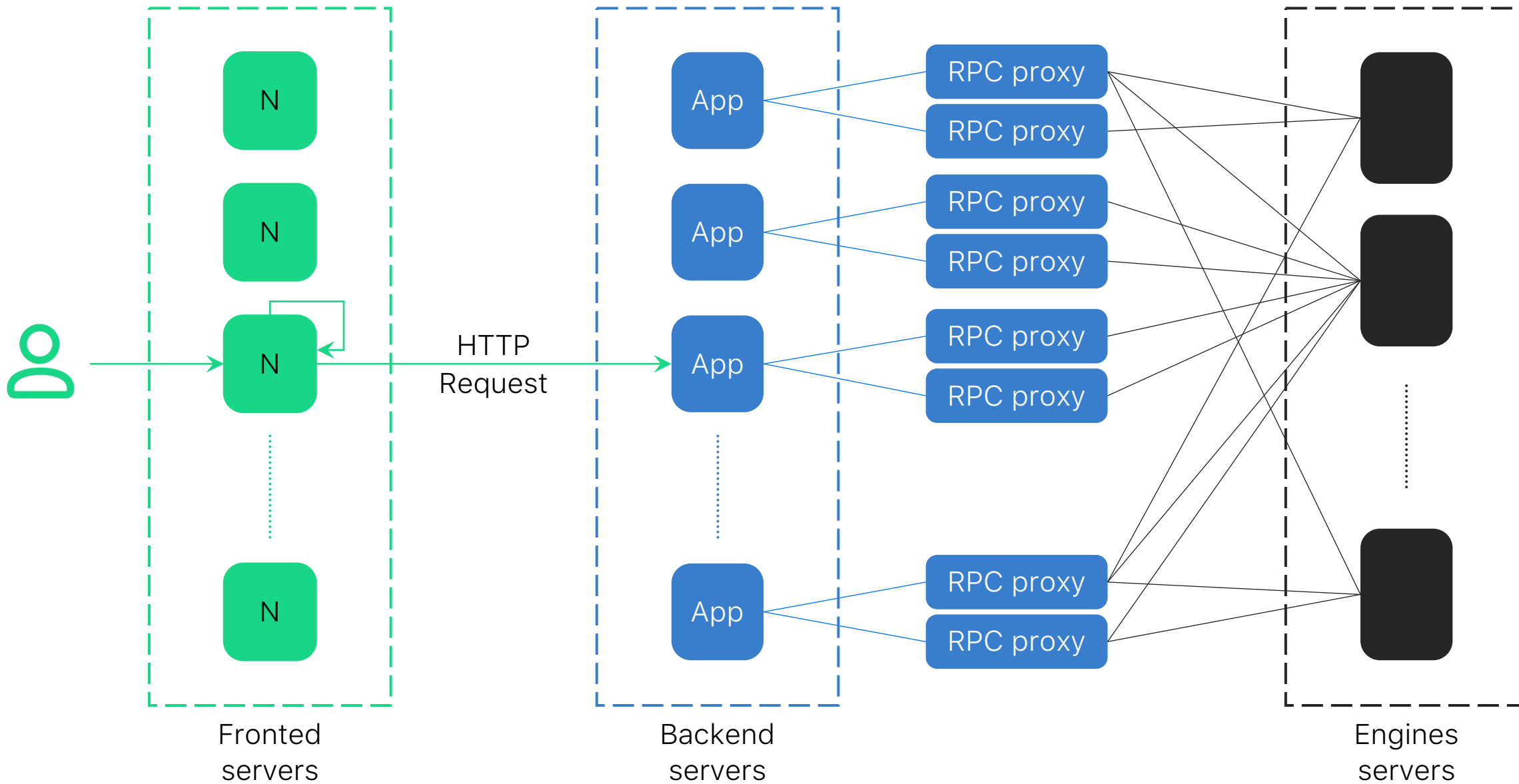
Поддержка продукта



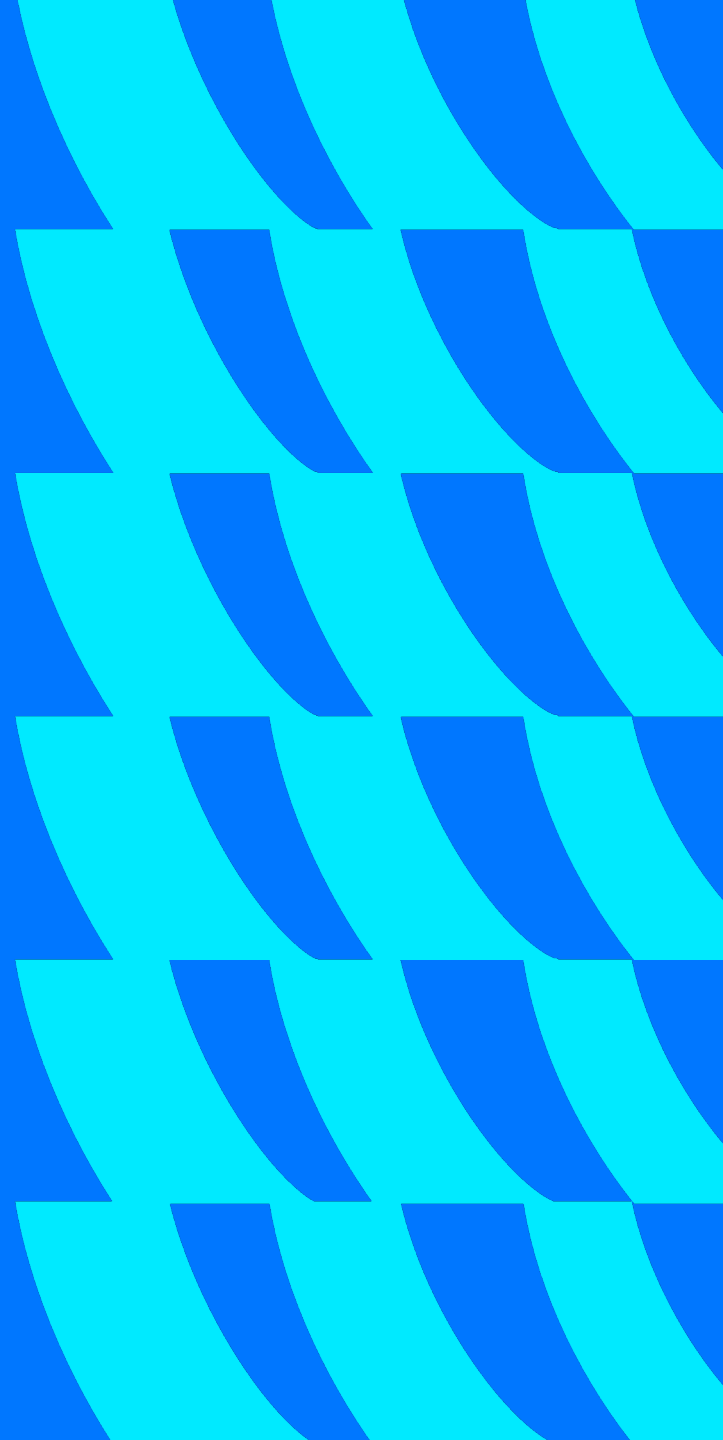
Утилизация ресурсов

Общая архитектура ВКонтакте





Движки ВКонтакте



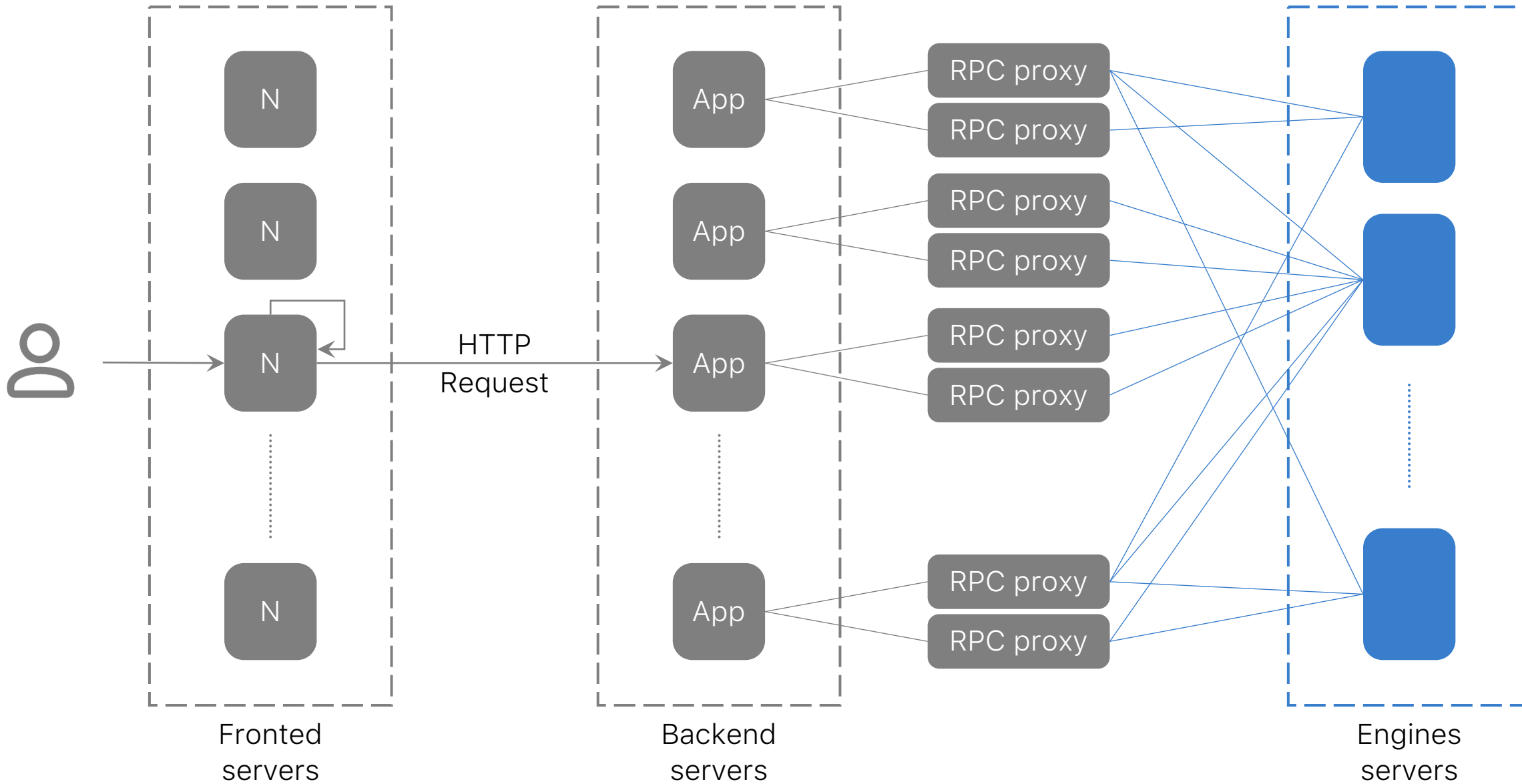
250 mpps

обрабатывают движки запросов в секунду

720 кластеров

кластер — отдельный этап движения

53 ДВИЖКА



Что делают движки?

1

Хранят и обрабатывают
все данные
пользователей

2

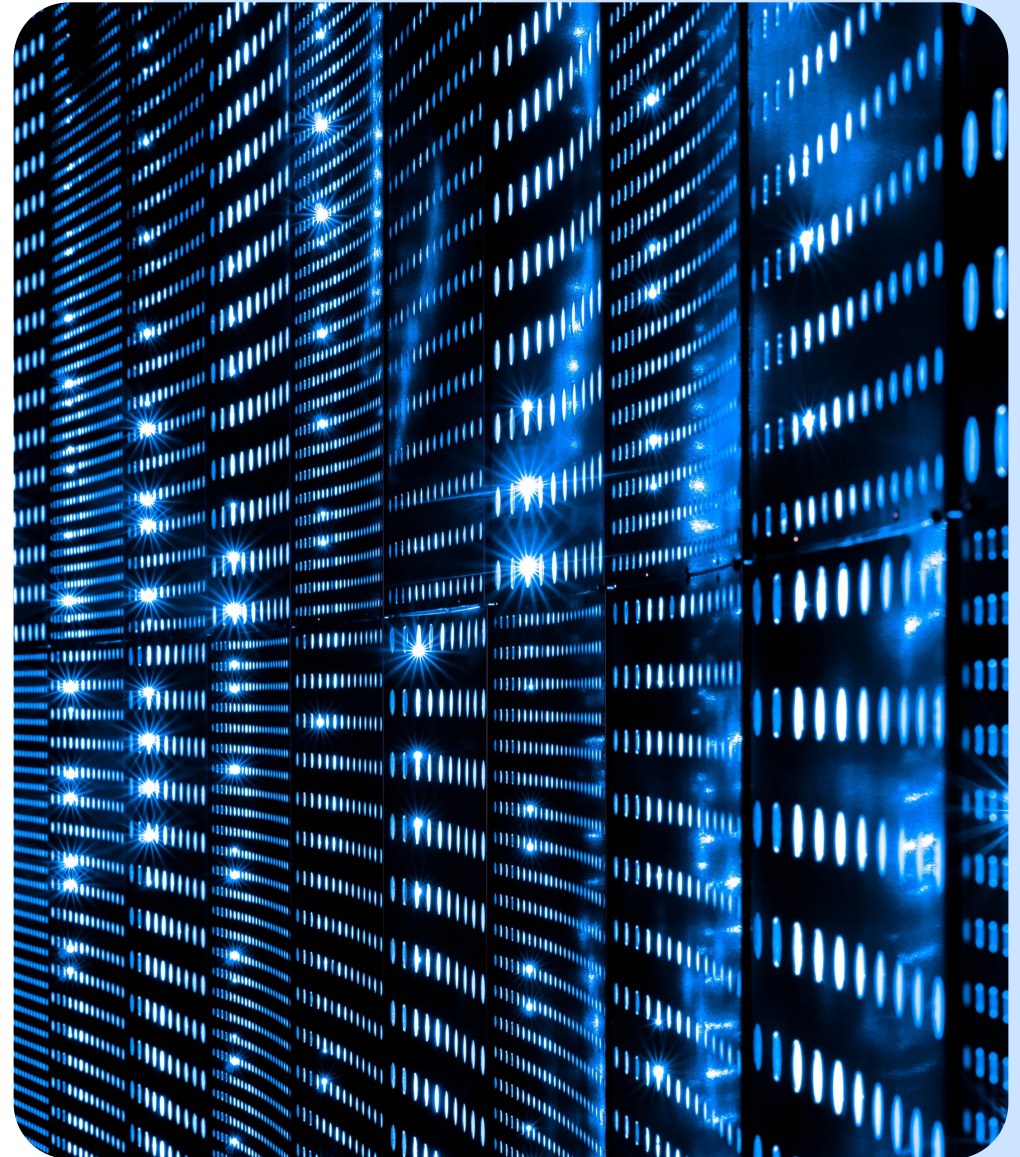
Контролируют
и реализуют
кеширующий слой

3

Хранят различные
служебные данные

4

Распространяют
конфигурацию



Почему мы

РАЗРАБАТЫВАЕМ

движки?

1

Гомогенность
эксплуатации

2

Экспертность
в разработке
и поддержке

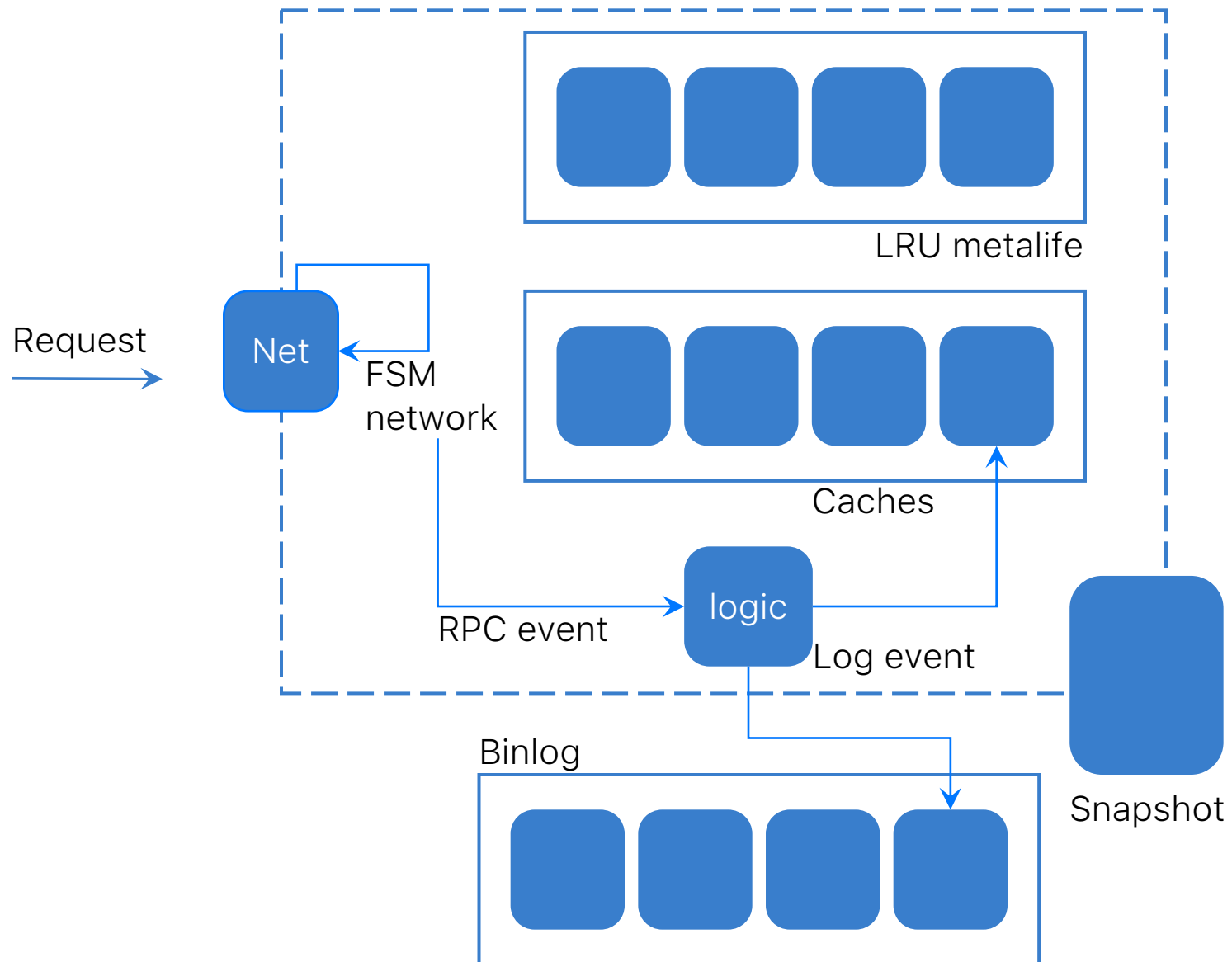
3

Знание
о хранимых
данных

4

Общая кодовая база
сетевого слоя
и работы с диском

ТИПОВОЙ ДВИЖОК



Binlog & snapshot

| OP | Users | Value |
|-----|-------|-------|
| add | 5015 | 1 |
| add | 5030 | 2 |
| add | 5015 | 2 |
| sub | 5030 | 2 |
| add | 5030 | 5 |

| Users | Value |
|-------|-------|
| 5015 | 3 |
| 5030 | 5 |

Indexation

| OP | Users | Value |
|-----|-------|-------|
| add | 5015 | 1 |
| add | 5030 | 2 |
| add | 5015 | 2 |
| sub | 5030 | 2 |
| add | 5030 | 5 |

Indexation



| Users | Value |
|-------|-------|
| 5015 | 3 |
| 5030 | 5 |



1.



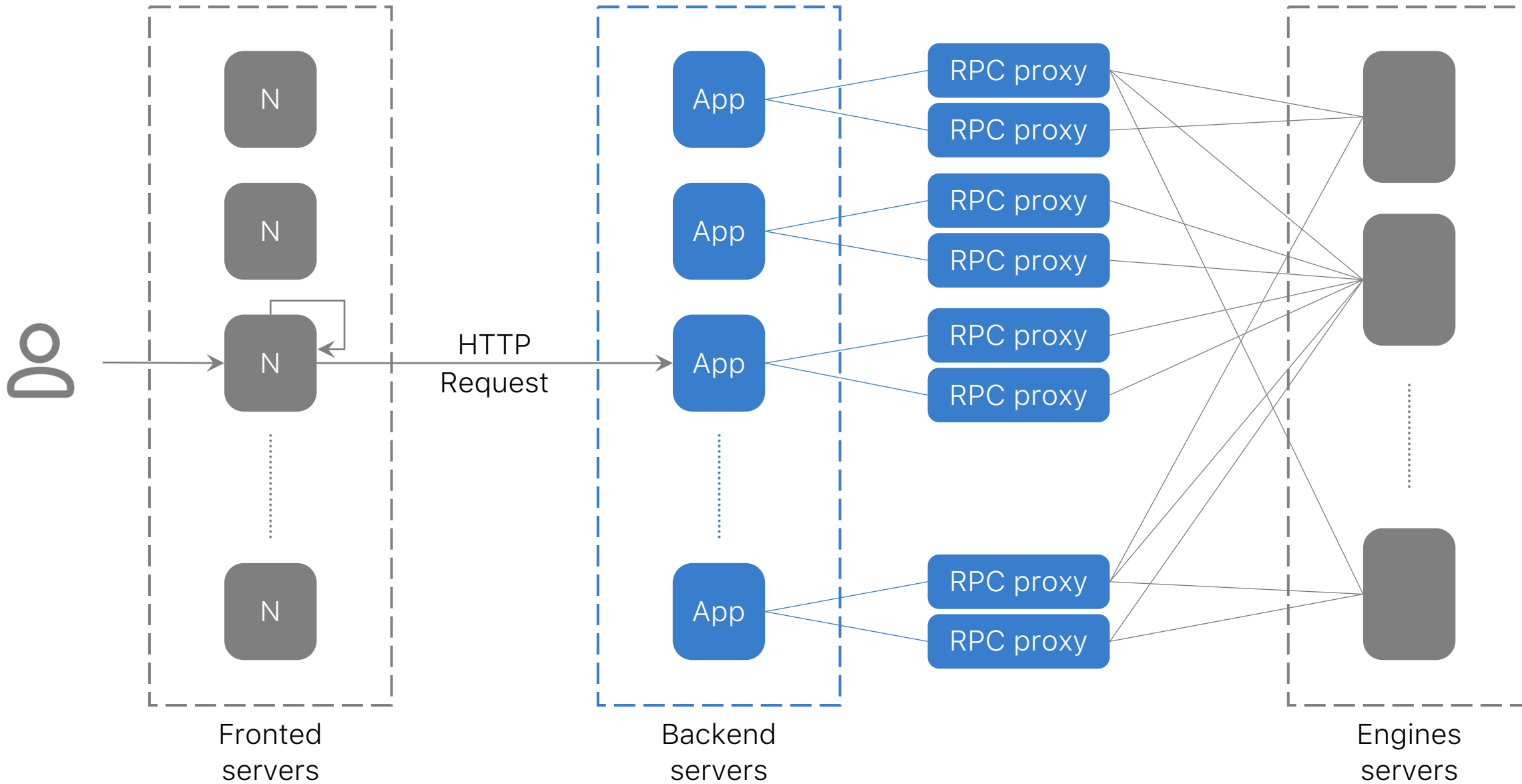
2.

Особенности ДВИЖКОВ

- 1 Однопоточные*
- 2 Данные и бизнес-логика рядом
- 3 Чтений сильно больше записей
- 4 Индексируются в отдельном процессе, используя COW операционной системы
- 5 Данные персистентные, но в основном IN-MEMORY
- 6 Скорость важнее консистентности на чтение

RPC-proxy

основной компонент для роутинга
и контроля распространения запросов



Топология связей

1

PRC-proxy знает
про каждый кластер

2

PRC-proxy
знает про шарды

3

PRC-proxy установлен
локально к каждому
APP-серверу

4

Занимает один
хоп по сети



ПРОБЛЕМА

mesh-топологии

1

При отказе одного из серверов общее время ответа кластера увеличивается

2

Начинает страдать время ответа API

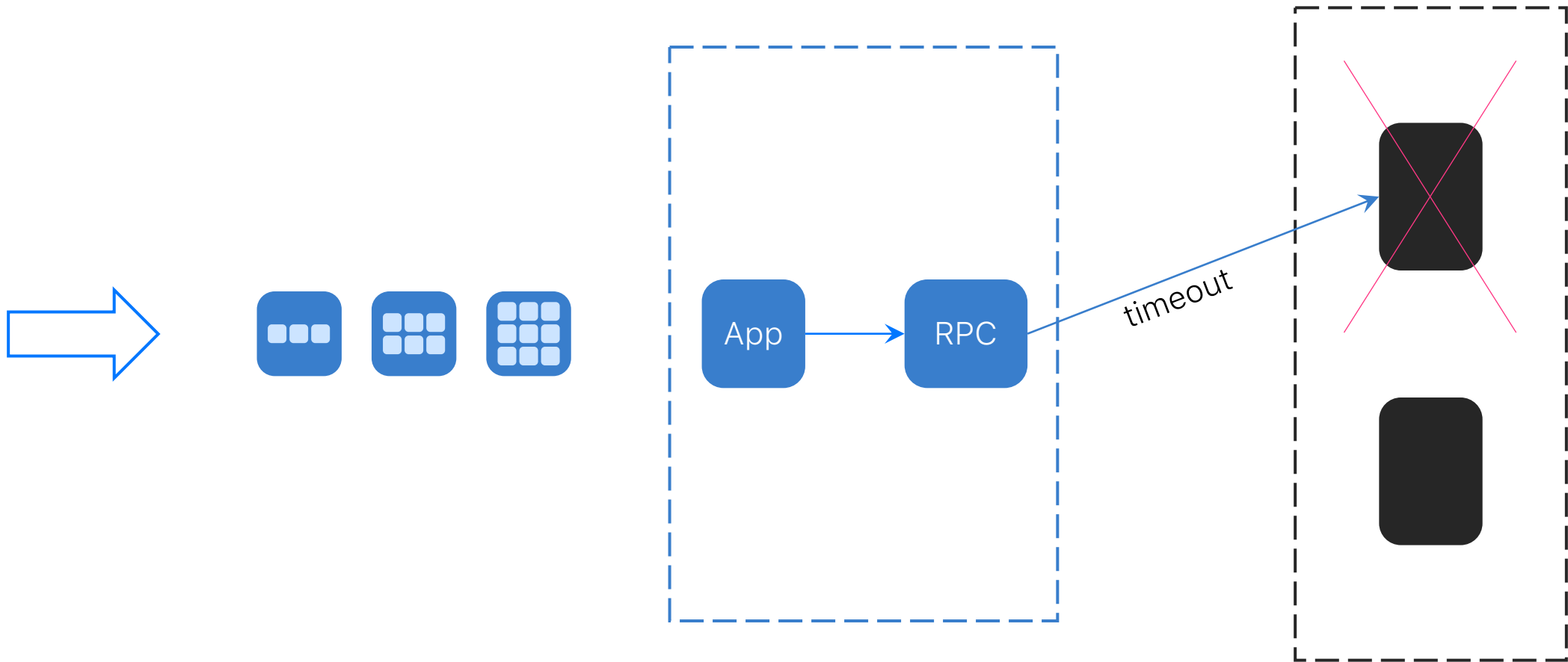
3

Заканчиваются ресурсы Backend-серверов

4

Получаем отказ в обслуживании всего сайта

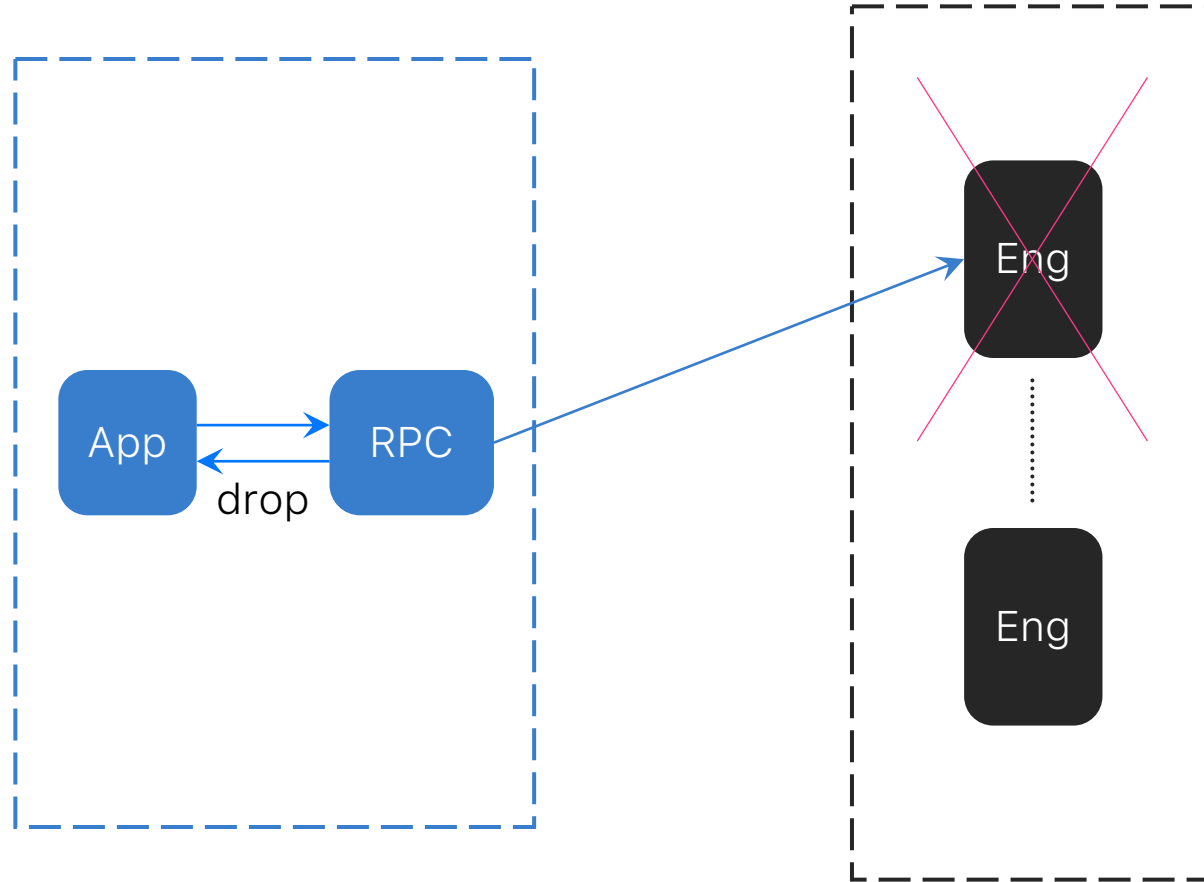


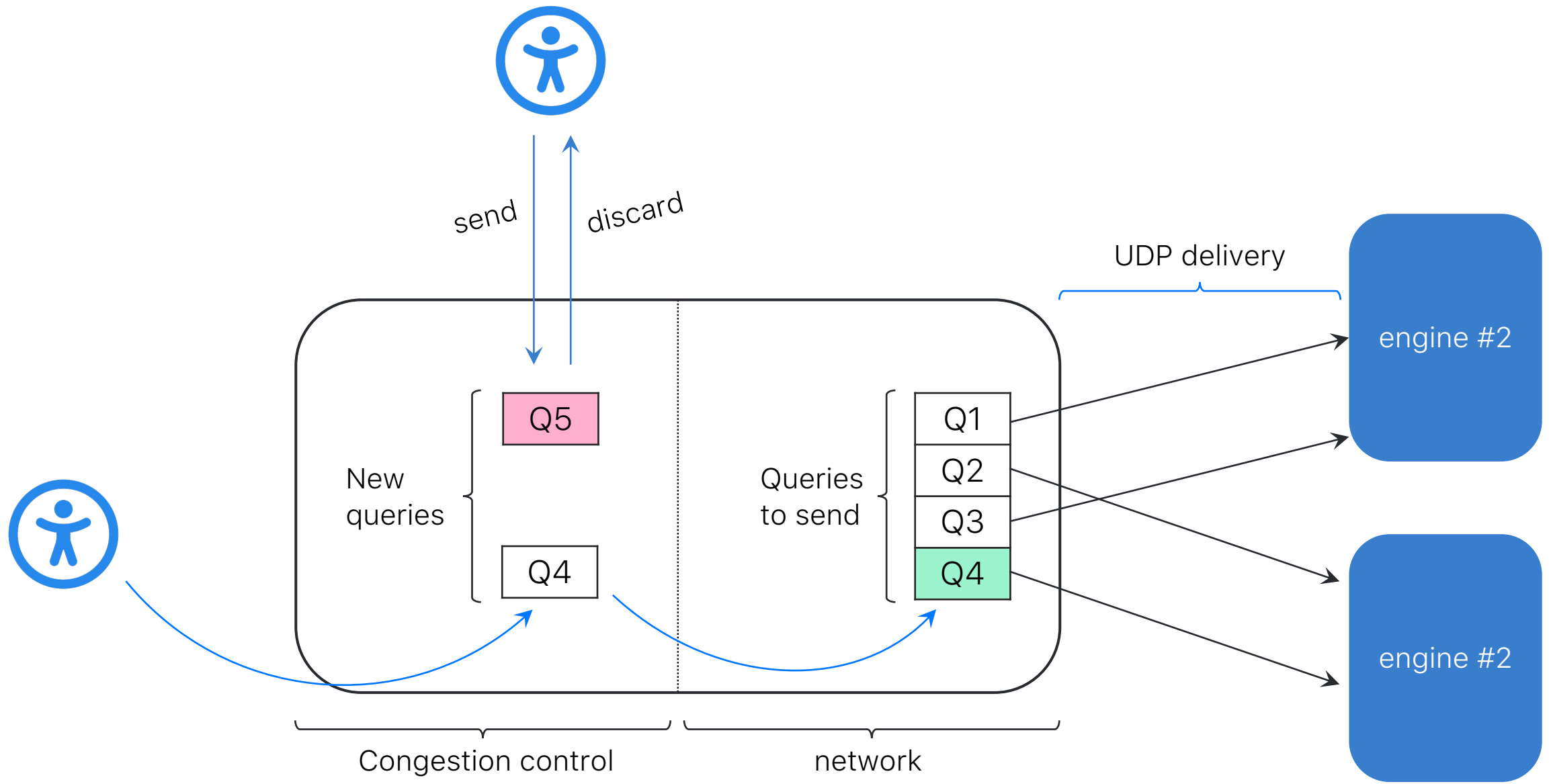


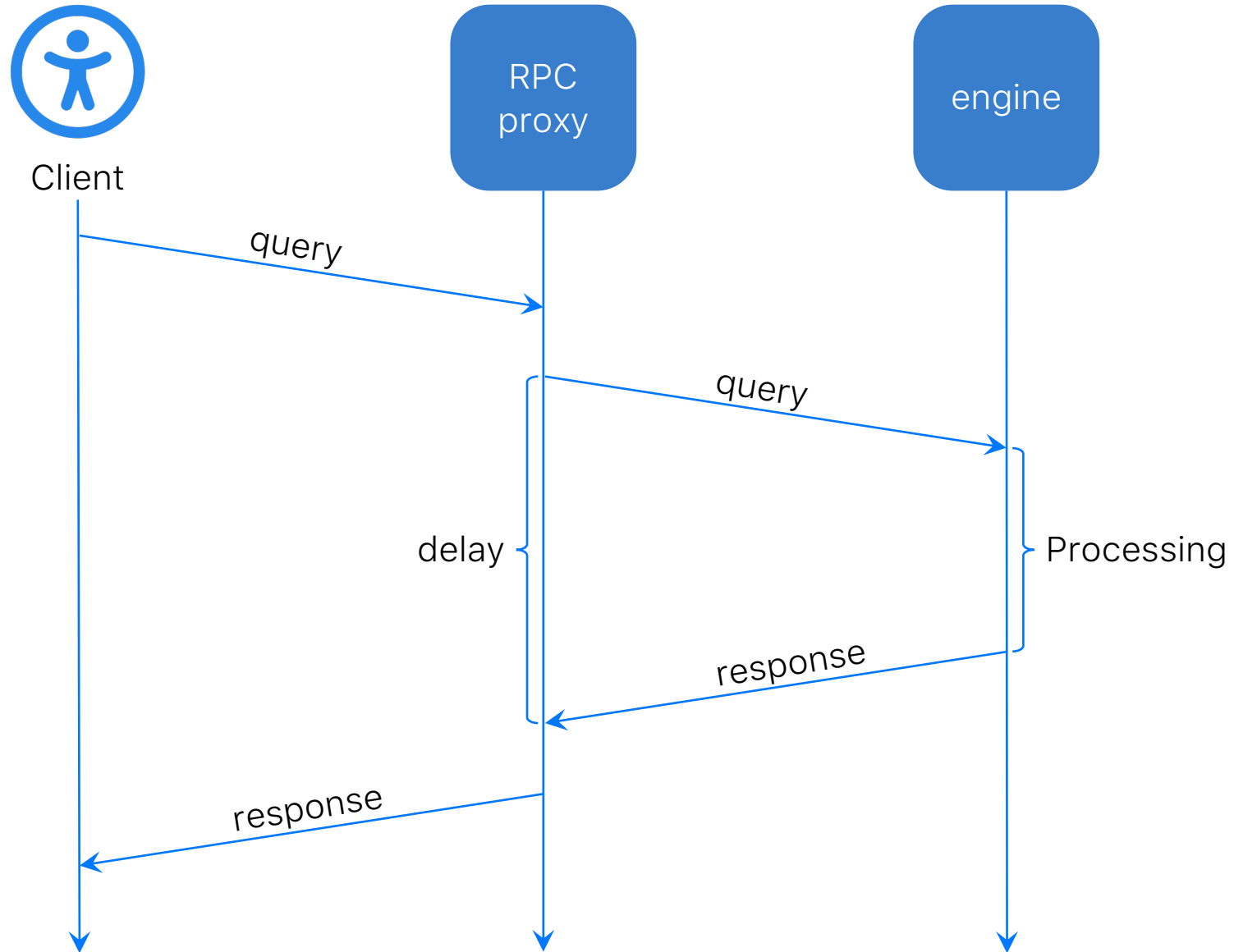
Решение — Congestion Control на application- уровне

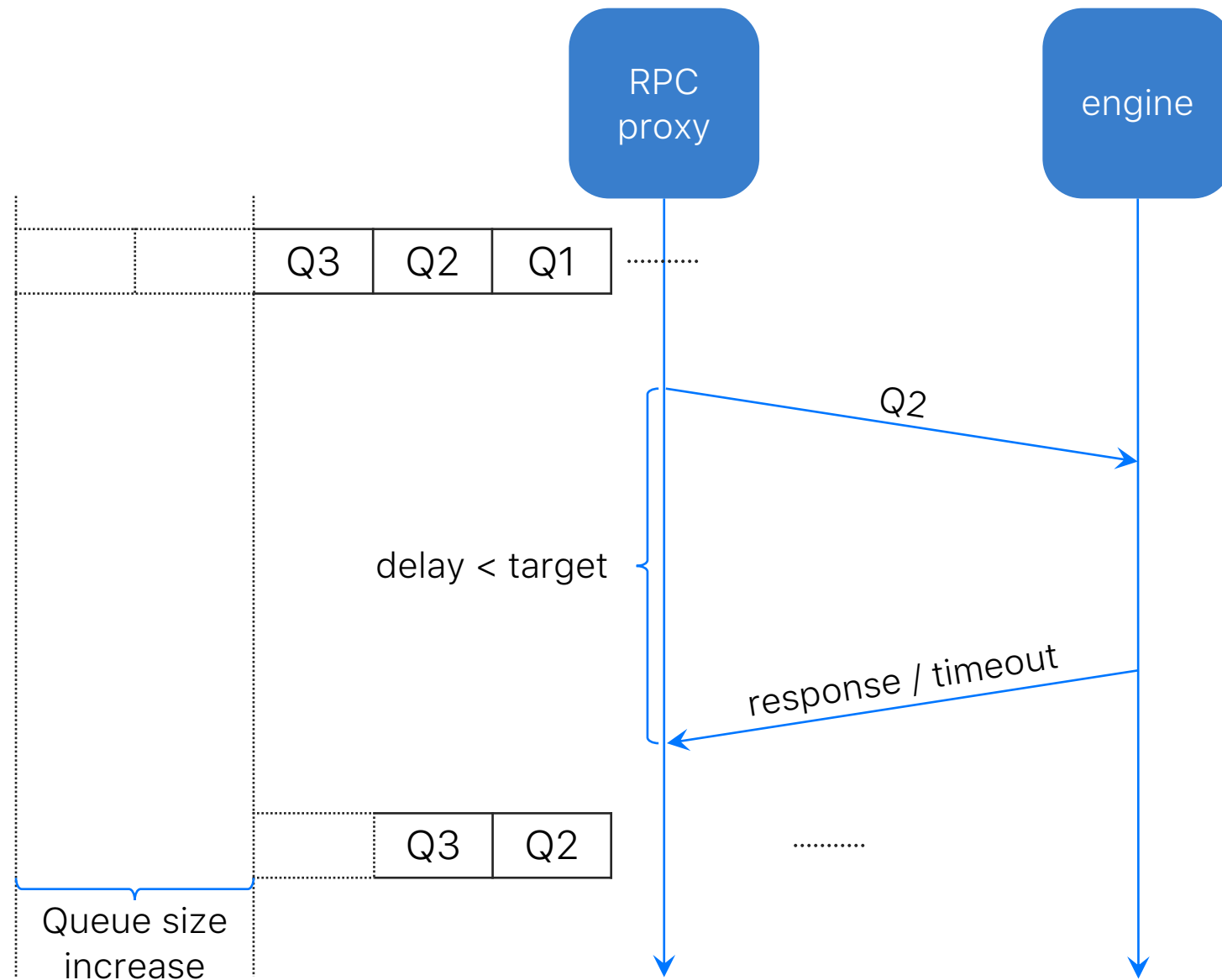
Идея: добавить на RPC-проxy
СС на уровне запросов

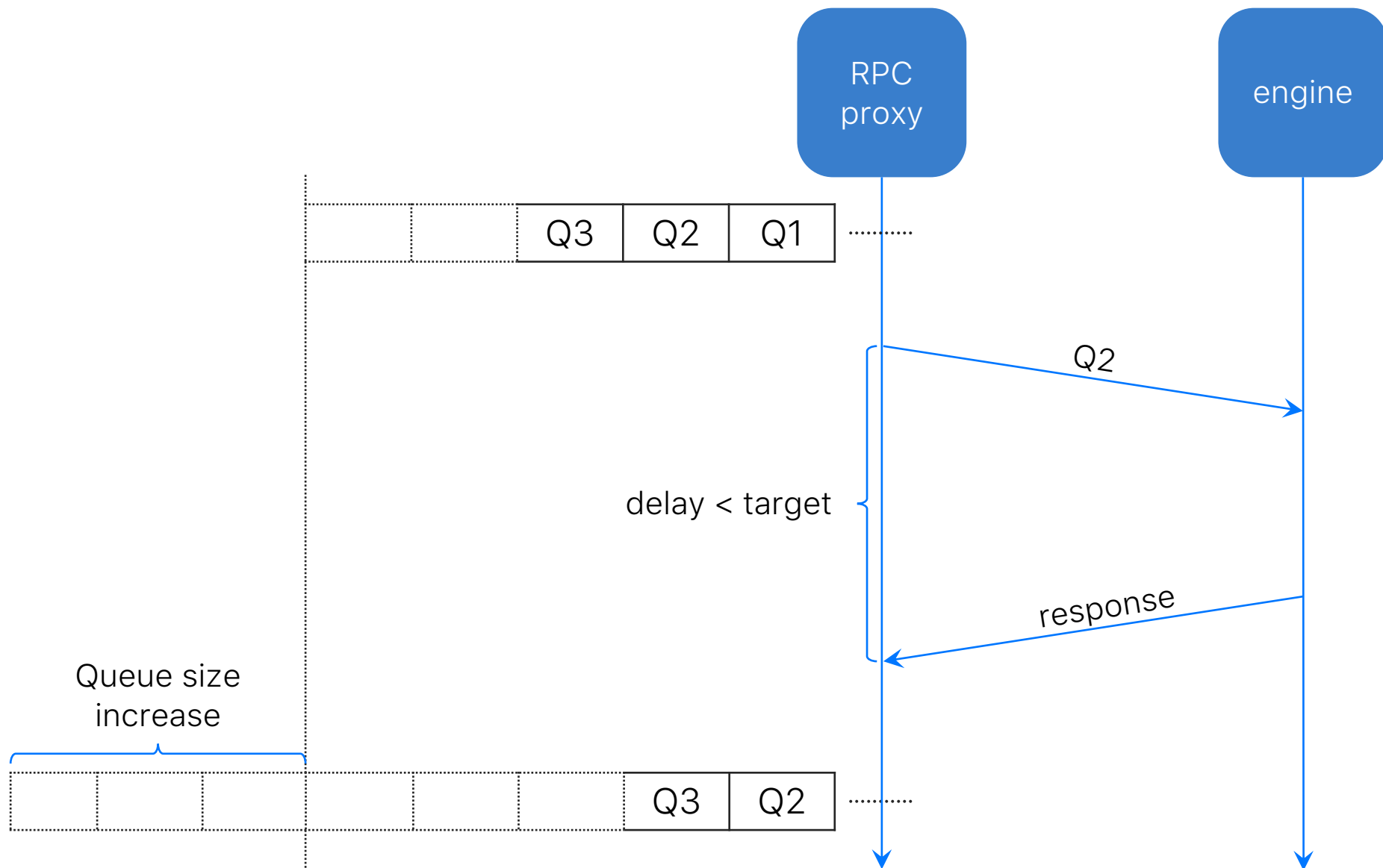
Если сервер начинает медленно
отвечать — в него начинает
приходить меньше запросов











discard

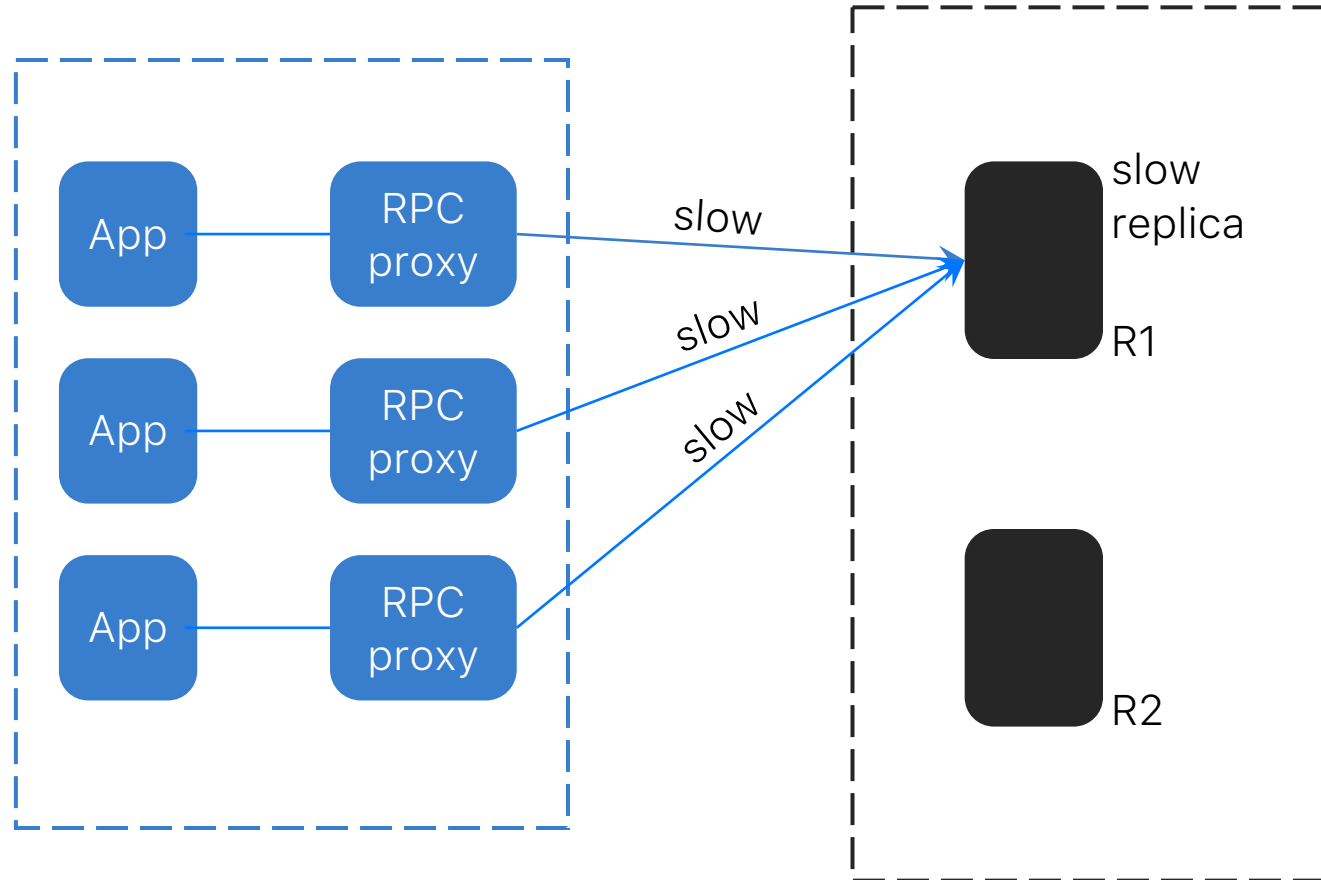


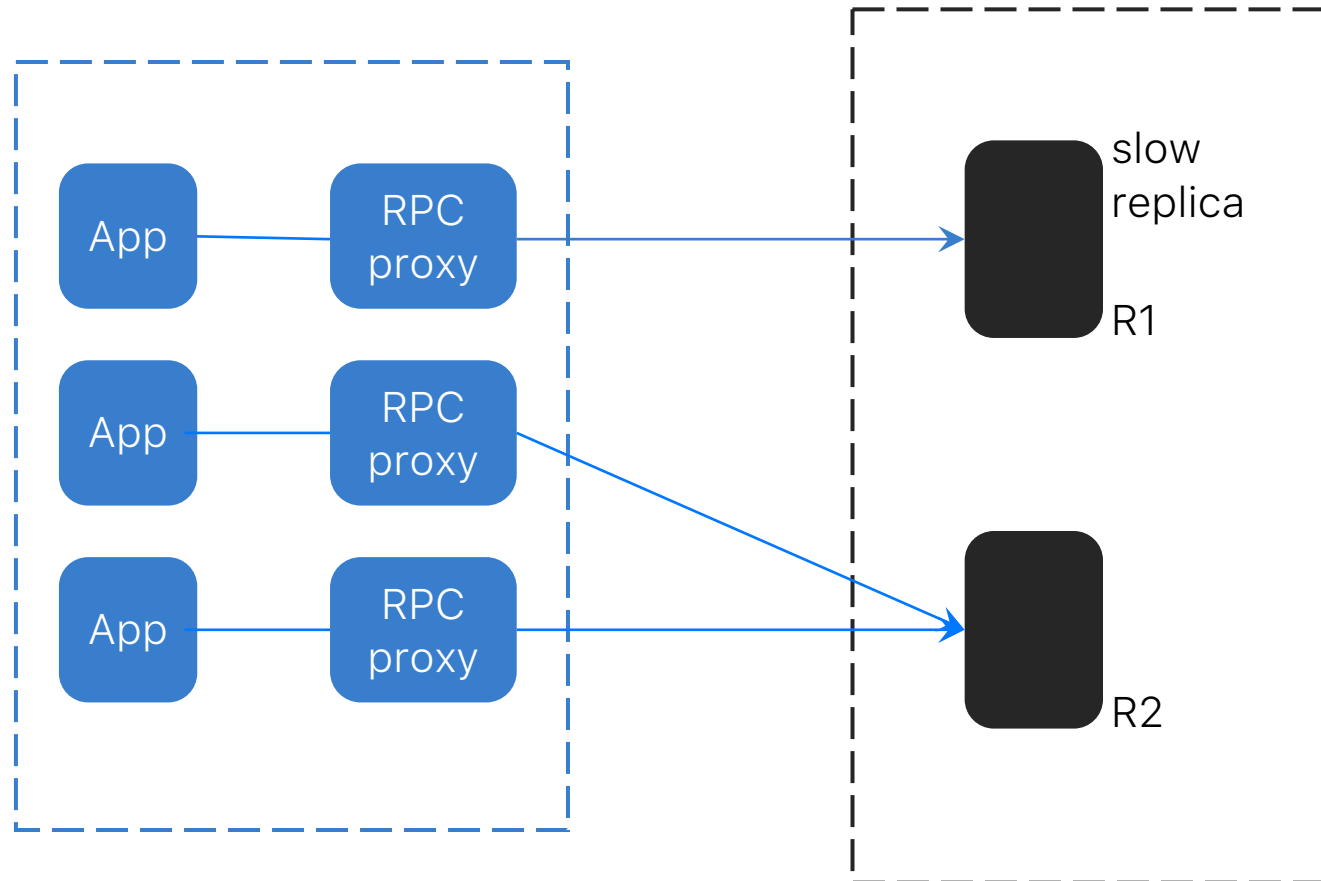
send



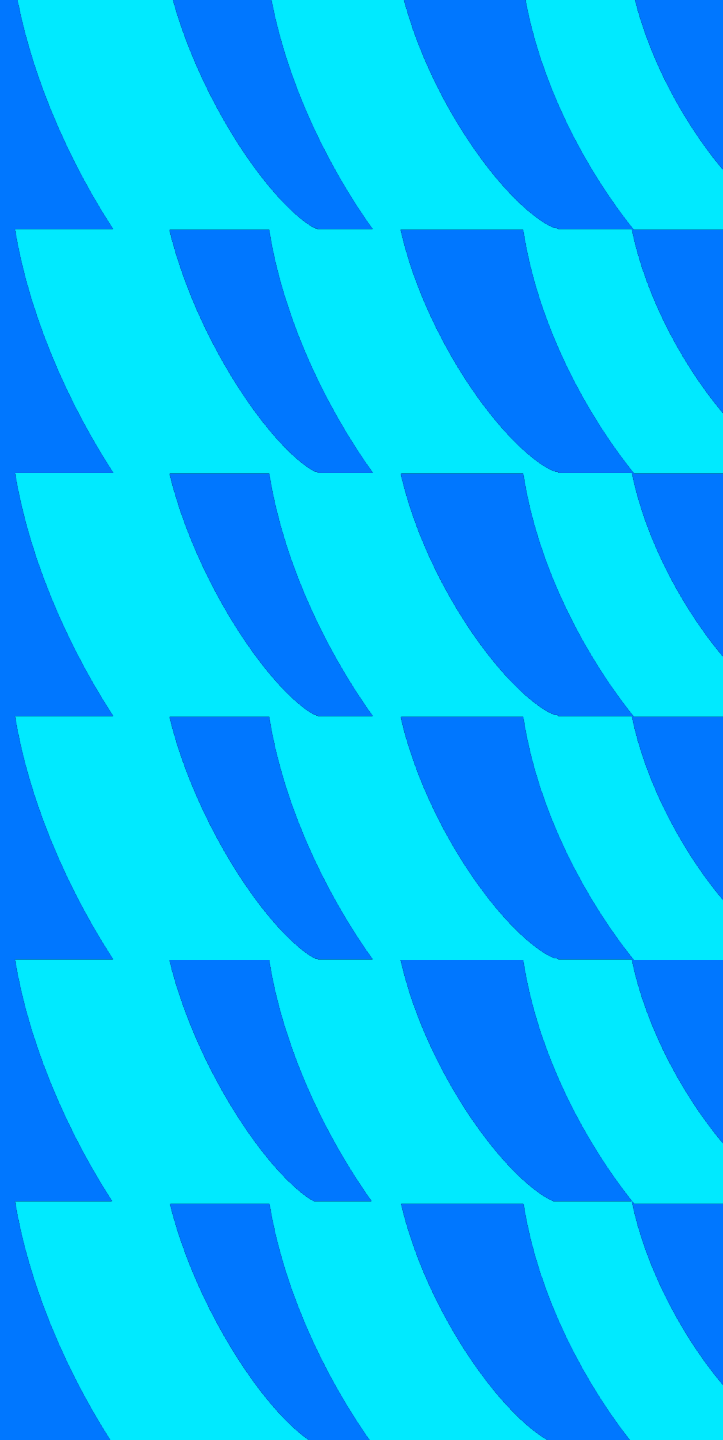
Application- level CC как балансир

Побочным эффектом является то, что механизмы CC позволяют уводить трафик с загруженной реплики





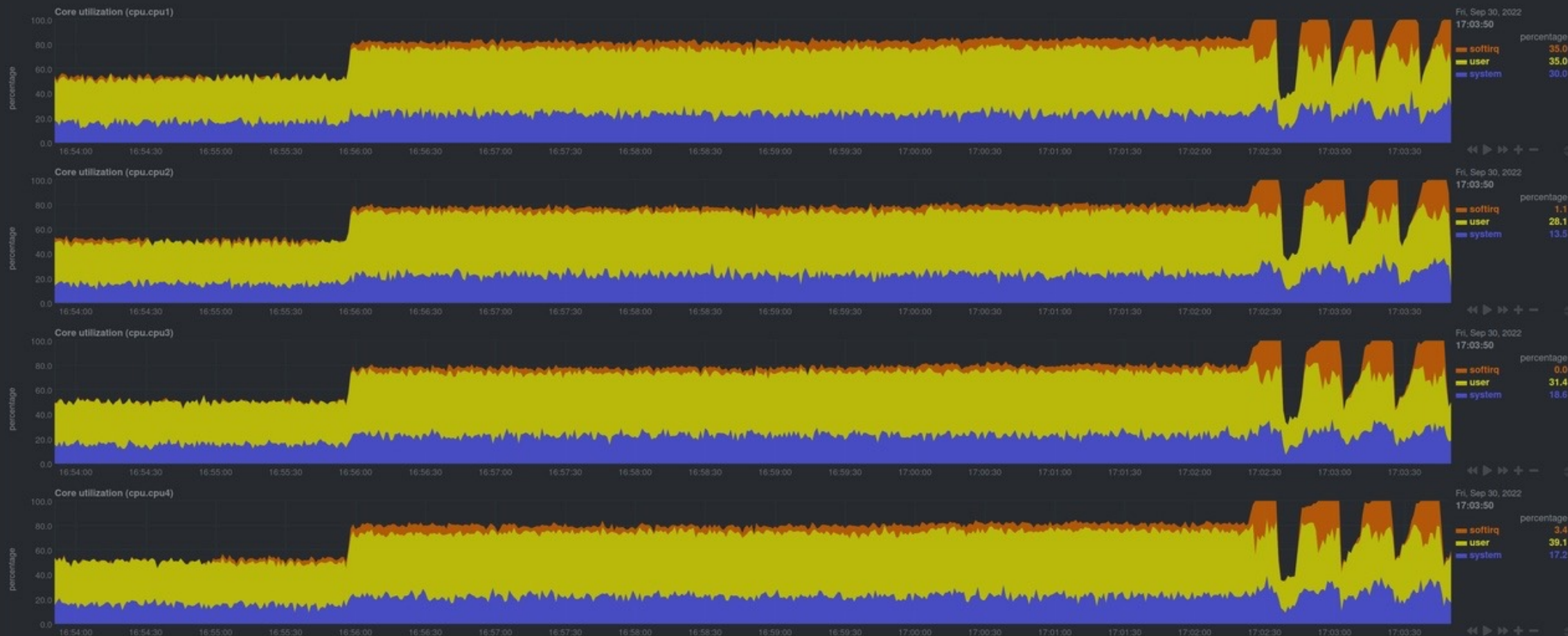
Результаты



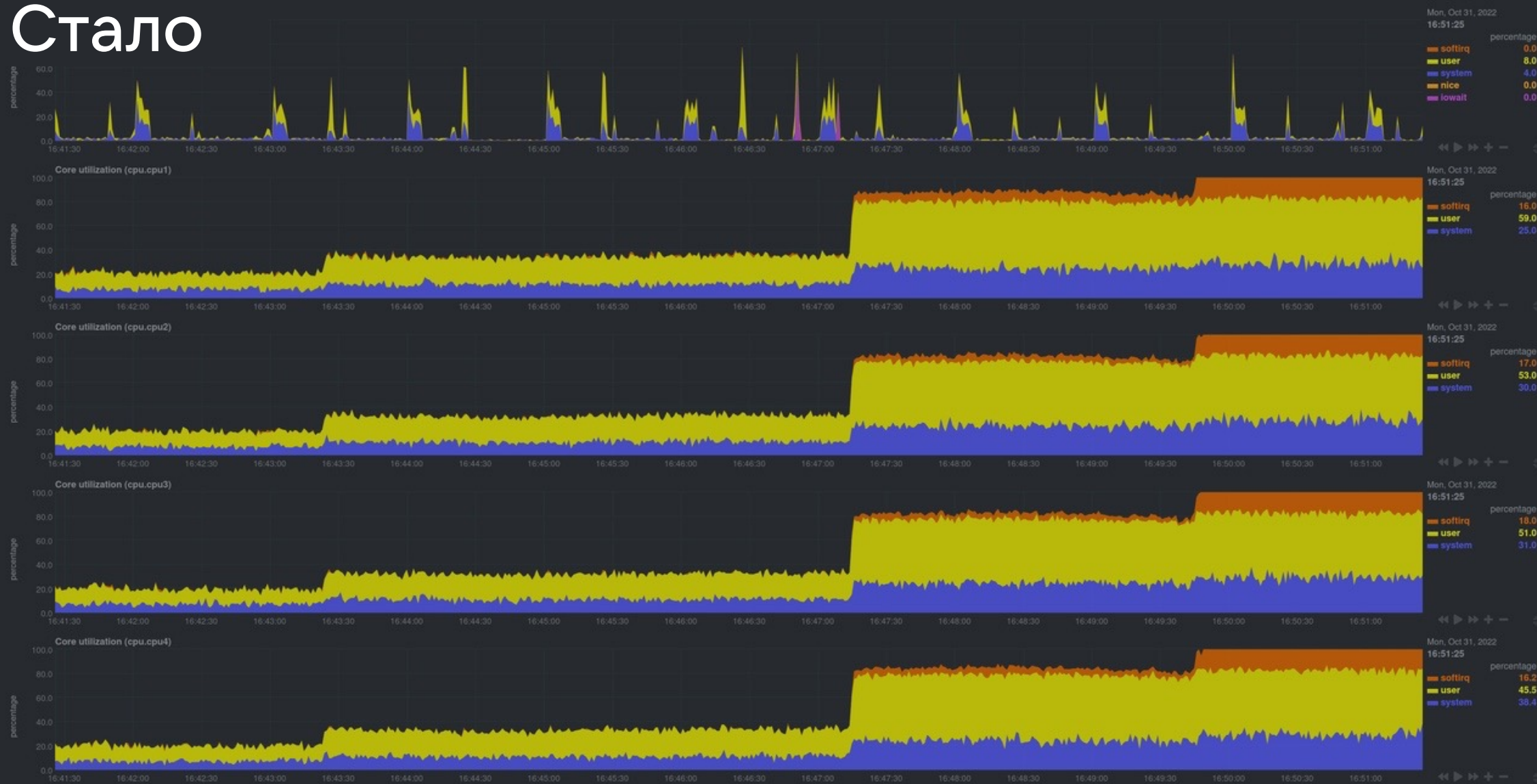
При нагрузке в 100% CPU
и утилизации канала сервер
отбивал часть запросов,
при этом не повышая
задержку успешных

При замедлении реплики
часть запросов идет
на соседнюю, снижая
при этом нагрузку
на загруженную

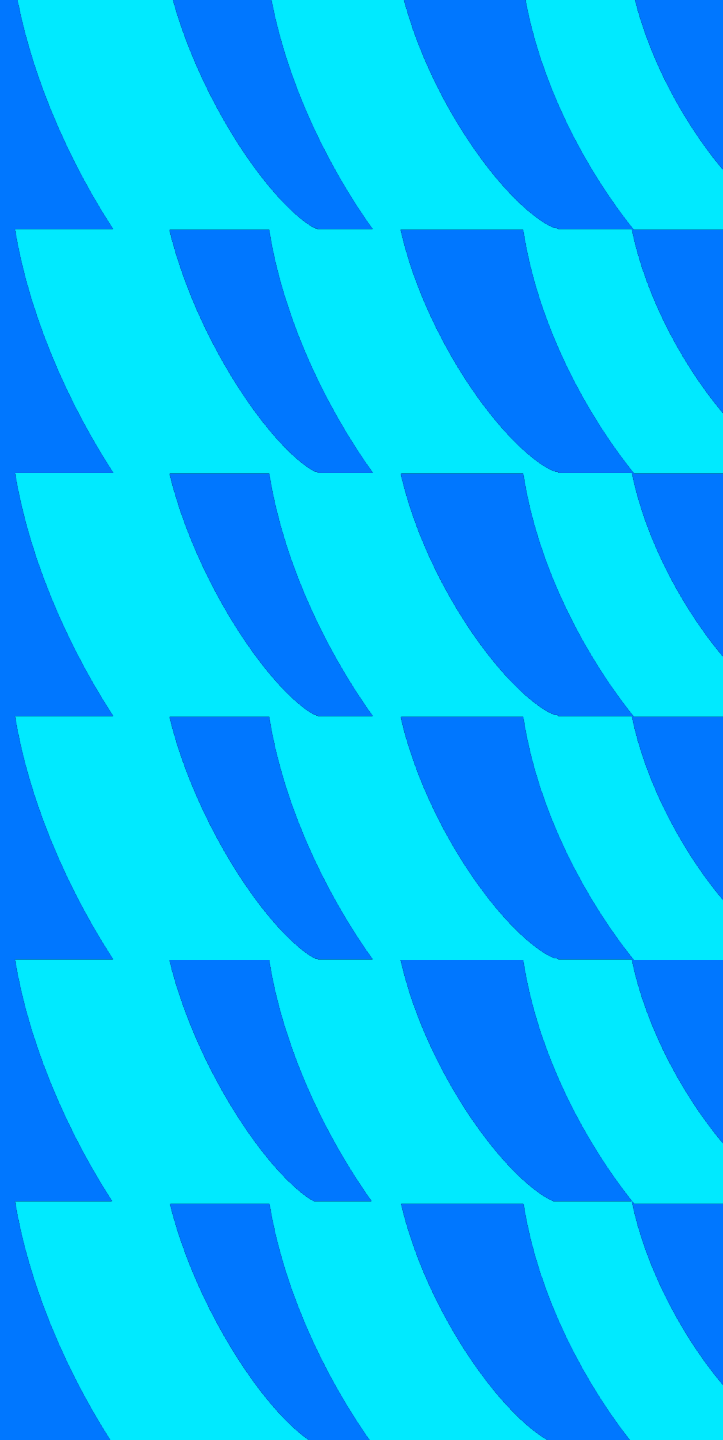
Было



Стало



Некоторые трюки с мемкешом





Memcached-shm

Порой ООМ-киллер или какая-то другая причина забирает лучших

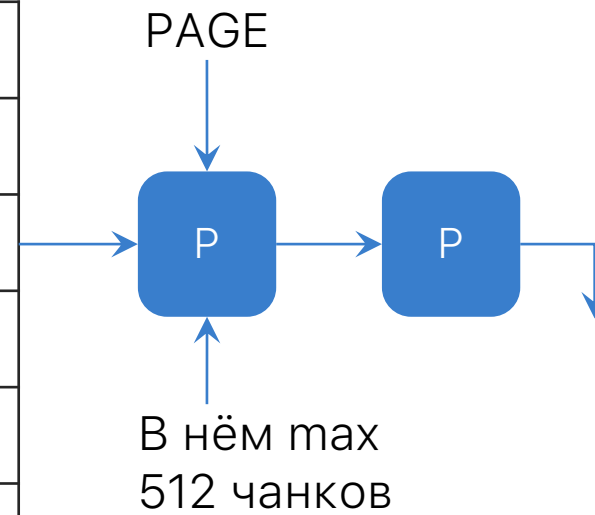
Для того чтобы не было так, был реализован простой аллокатор, использующий shared-memory

Memcached-shm: allocator

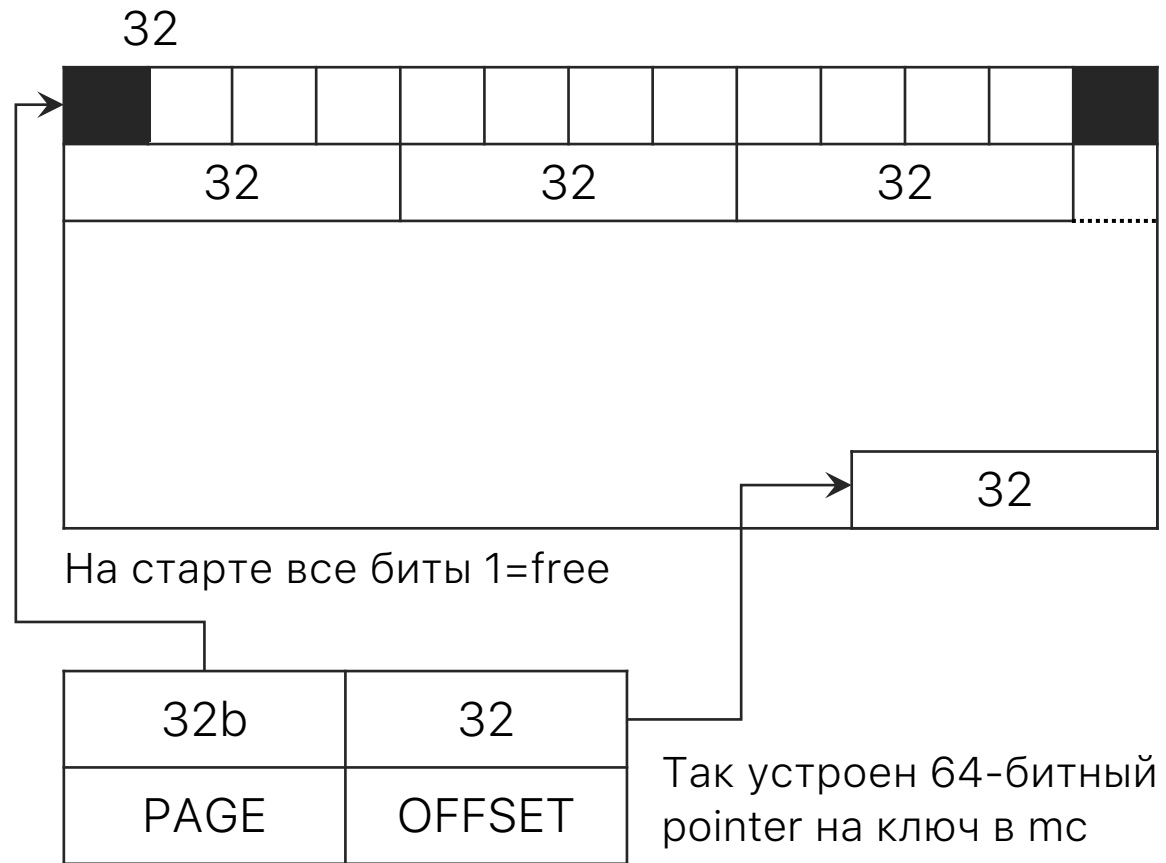
Таблица sizeclass-08:

| | |
|----|-----|
| 8 | PTR |
| 16 | PTR |
| 24 | PTR |
| 32 | |
| | |
| | |

Binsearch
в момент аллокации

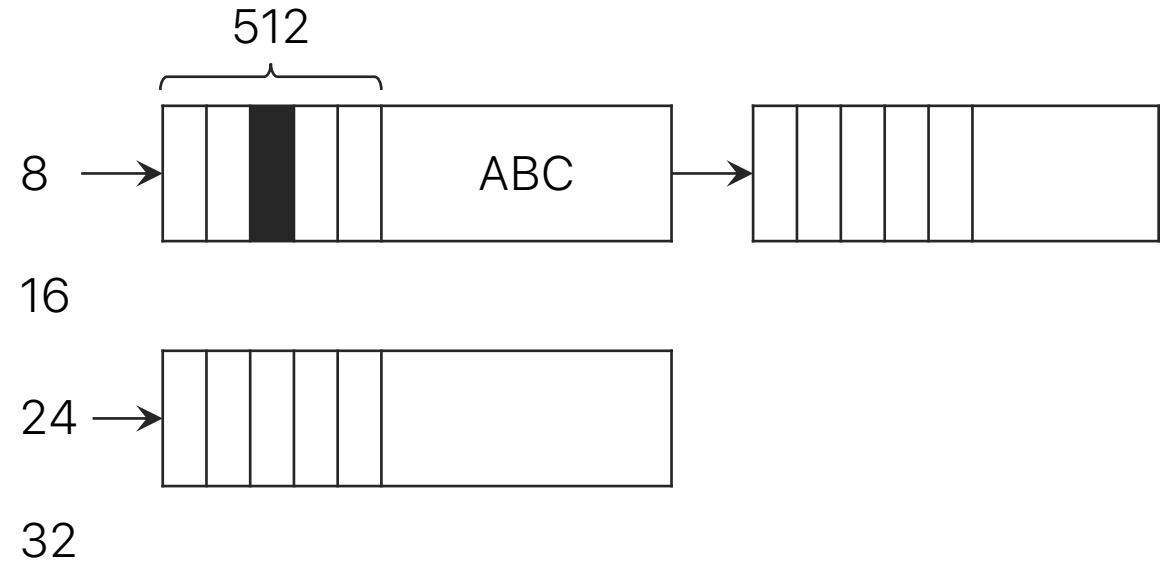


Memcached-shm: allocator



Так устроен 64-битный pointer на ключ в ms

Memcached-shm: allocator



Memcached-shm: статистика аллокаций

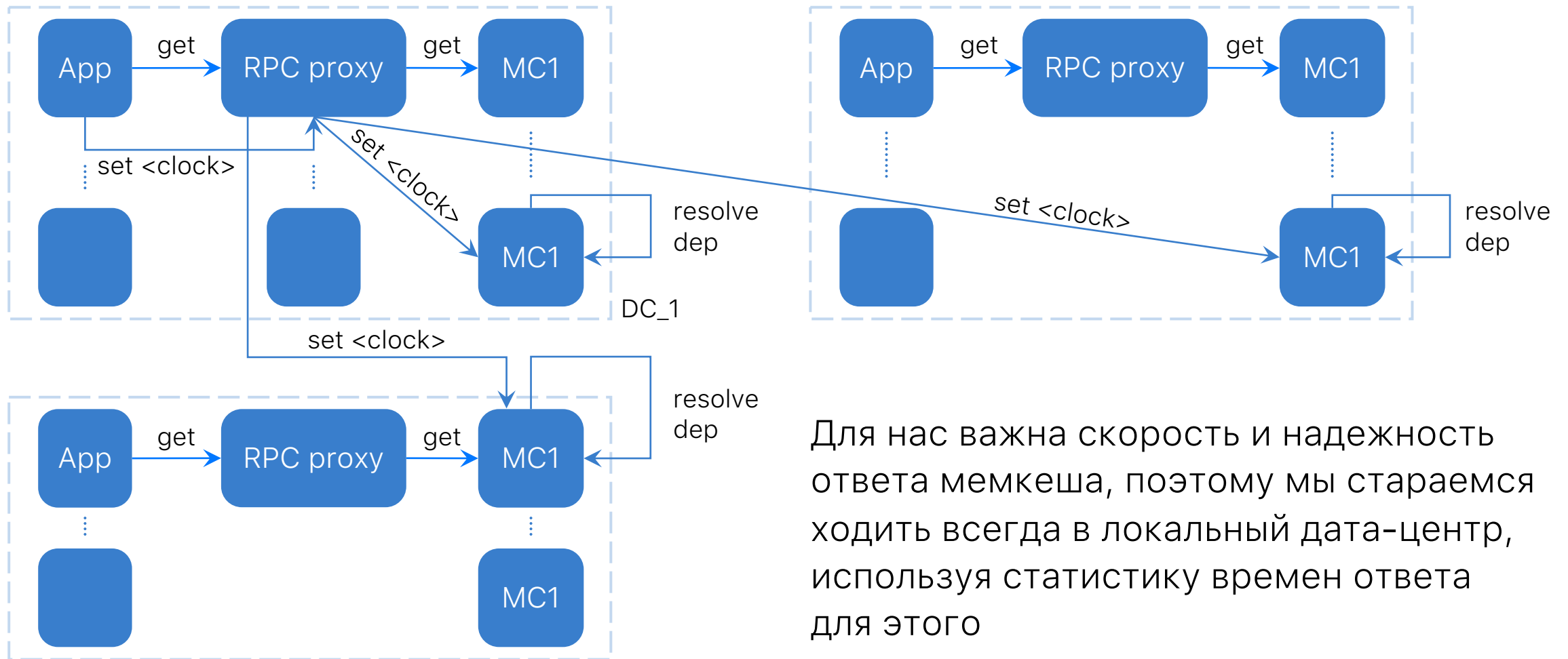
| sizeclass | lru_size_evaluation | lru_updates | lru_drops | lru_drops_ffff |
|-----------|---------------------|-------------|-----------|----------------|
| 0 | 0 | 0 | 0 | 0 |
| 1 | 6107 | 504133 | 0 | 0 |
| 2 | 82140 | 13986554 | 11698 | 11698 |
| 3 | 1788281 | 94169784 | 483222 | 483222 |
| 4 | 2849519 | 53064428 | 1551886 | 748672 |
| 5 | 643132 | 24328205 | 871995 | 163772 |
| 6 | 249453 | 23920864 | 284862 | 58490 |
| 7 | 456909 | 12624350 | 340068 | 116980 |
| 8 | 343046 | 28121403 | 567001 | 81886 |
| 9 | 53230 | 4440532 | 56235 | 11698 |
| 10 | 74278 | 4987243 | 20728 | 11698 |
| 11 | 88363 | 11146533 | 43233 | 17663 |
| 12 | 48847 | 10969569 | 35208 | 11698 |
| 13 | 88861 | 11663046 | 36099 | 23384 |
| 14 | 44527 | 7182312 | 12624 | 11698 |
| 15 | 59629 | 17415728 | 12251 | 11698 |
| 16 | 38810 | 6643563 | 21207 | 0 |

How I feel
when
I share
anything



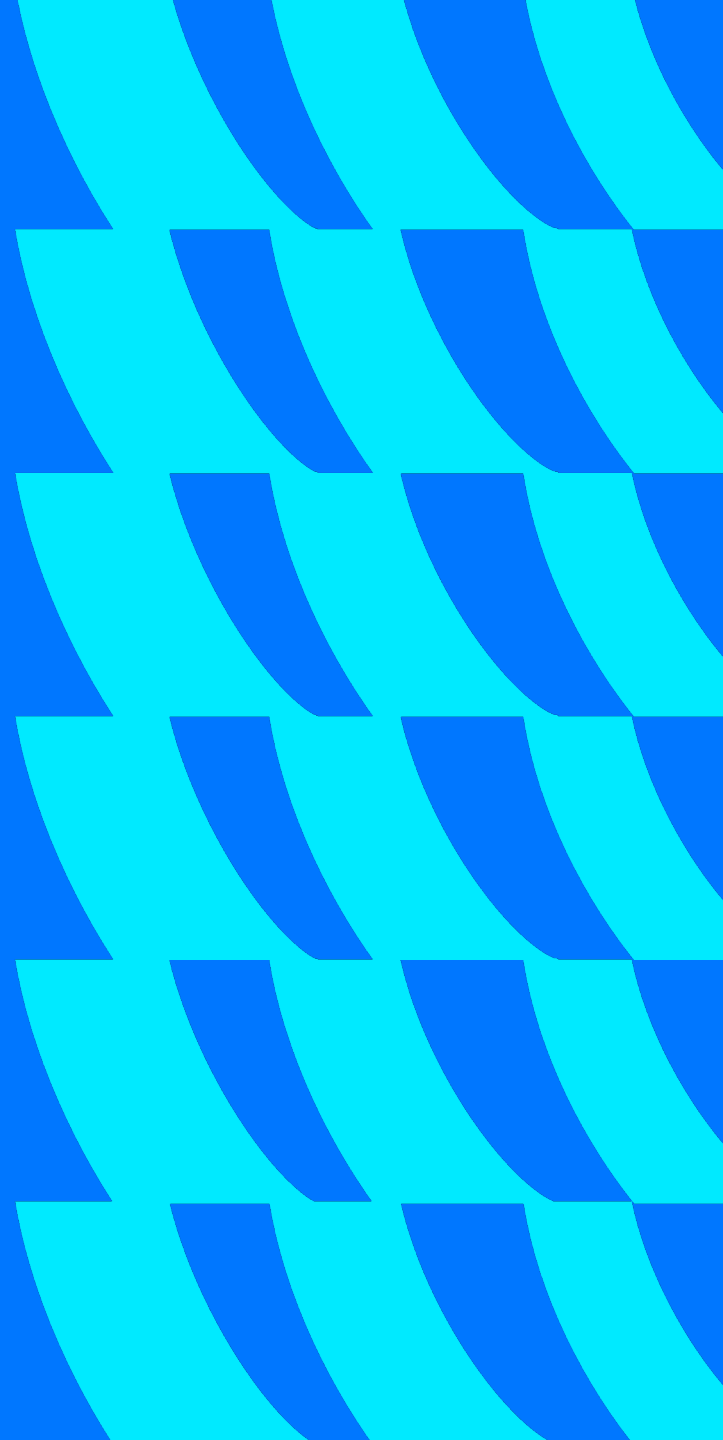
I serve the Soviet Union

Memcached-replication

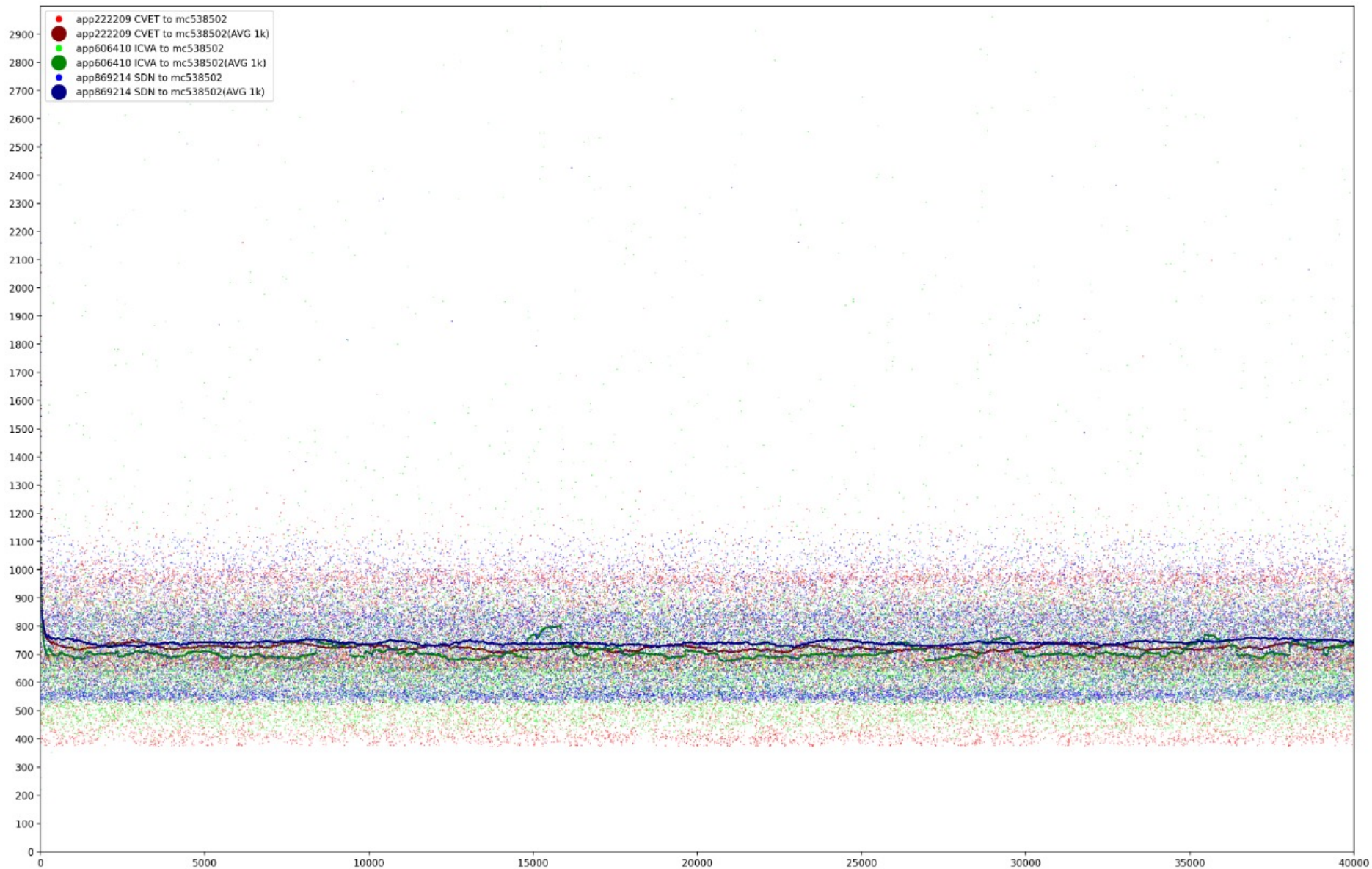


Для нас важна скорость и надежность ответа мемкеша, поэтому мы стараемся ходить всегда в локальный дата-центр, используя статистику времен ответа для этого

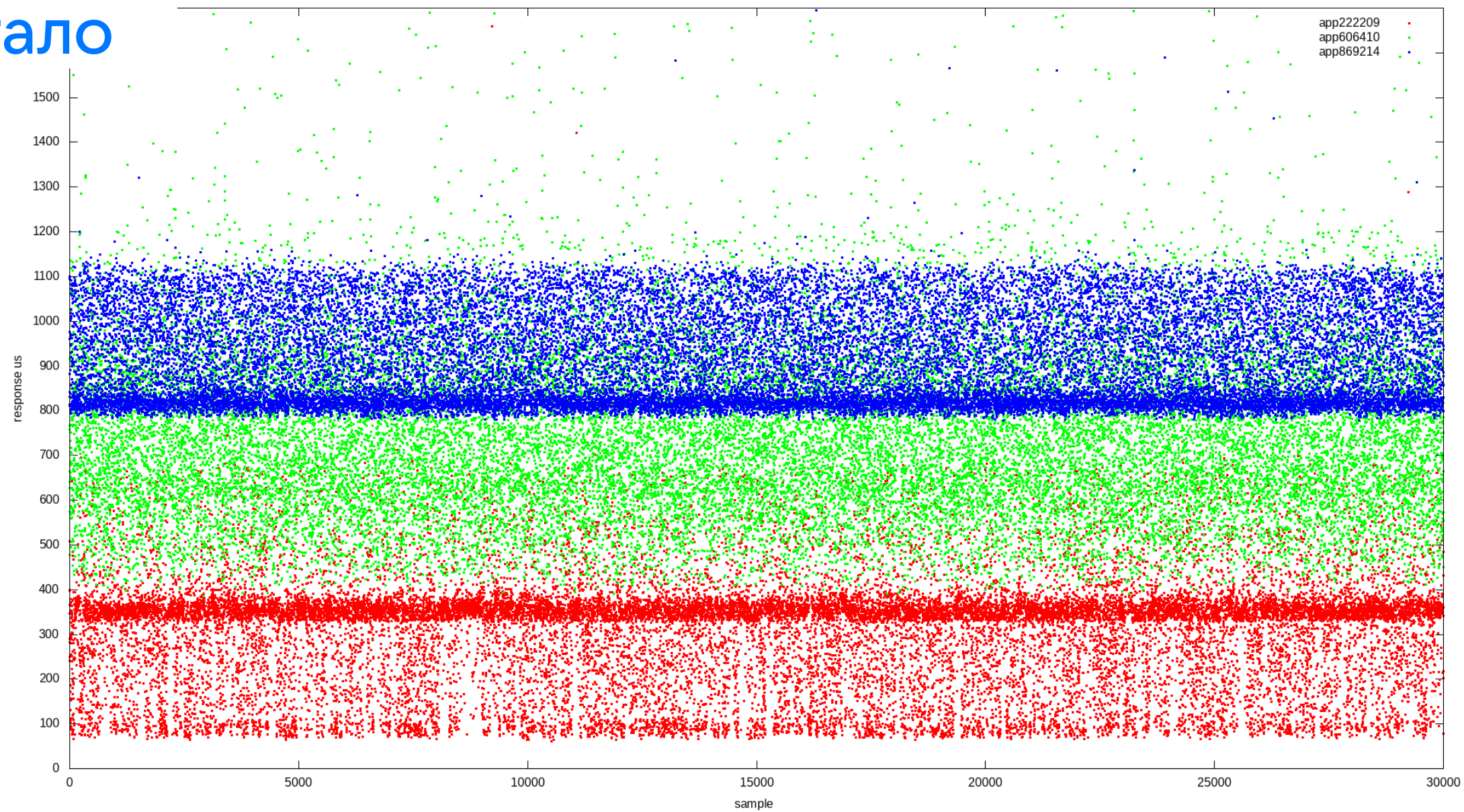
Почему
ms-replicated –
это важно?



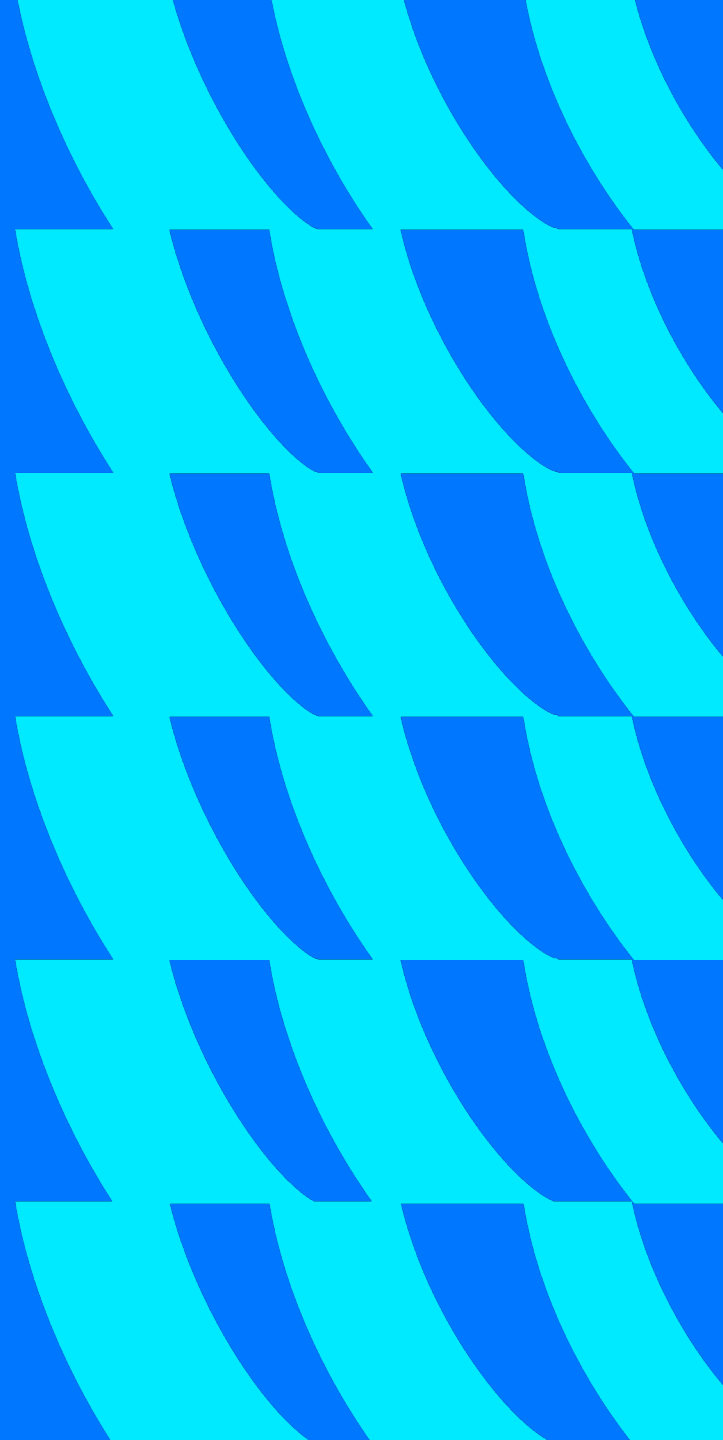
Было



Стало



Бывало
разное





РВАНЁТ

НЕ
ДОЛЖНО...

When you watch the video
at 2x speed to save time
but it is still $O(n)$



Нарвались на `list::size()`
за $O(n)$ в продакшне



Проблема

Часть кода собиралась старым компилятором, где несмотря на 11-ый стандарт, функция взятия размера выполнялась за линию

Что не преминуло сказаться на графиках



Проблема

Часть кода собиралась старым компилятором, где несмотря на 11-ый стандарт, функция взятия размера выполнялась за линию

Что не преминуло сказаться на графиках



Решение

Сейчас перешли на 20-ый стандарт и новые компиляторы, а для старых сборок где нужно подкладываем сборную солянку из старого glibc

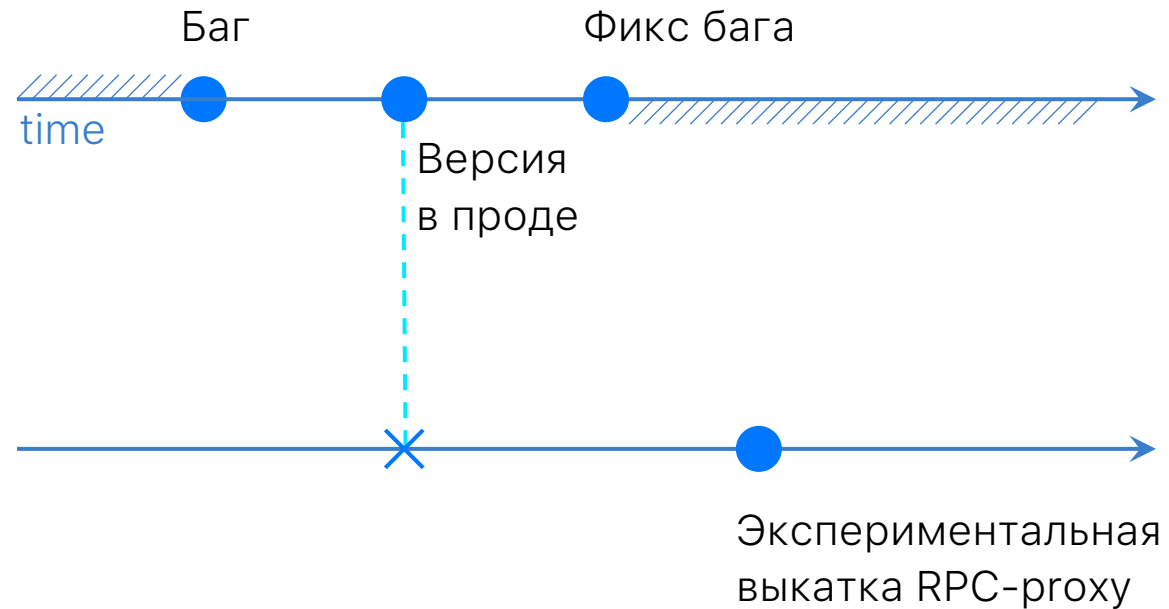
Да, контейнеров у нас для движков нет



«Пакет смерти»

Самолет не падает при отказе одного двигателя

Так и здесь – при определенной конфигурации и определенной версии нашелся пакет, который вынес нам кеширующий слой



«Пакет смерти»: решение



Фазировали сеть

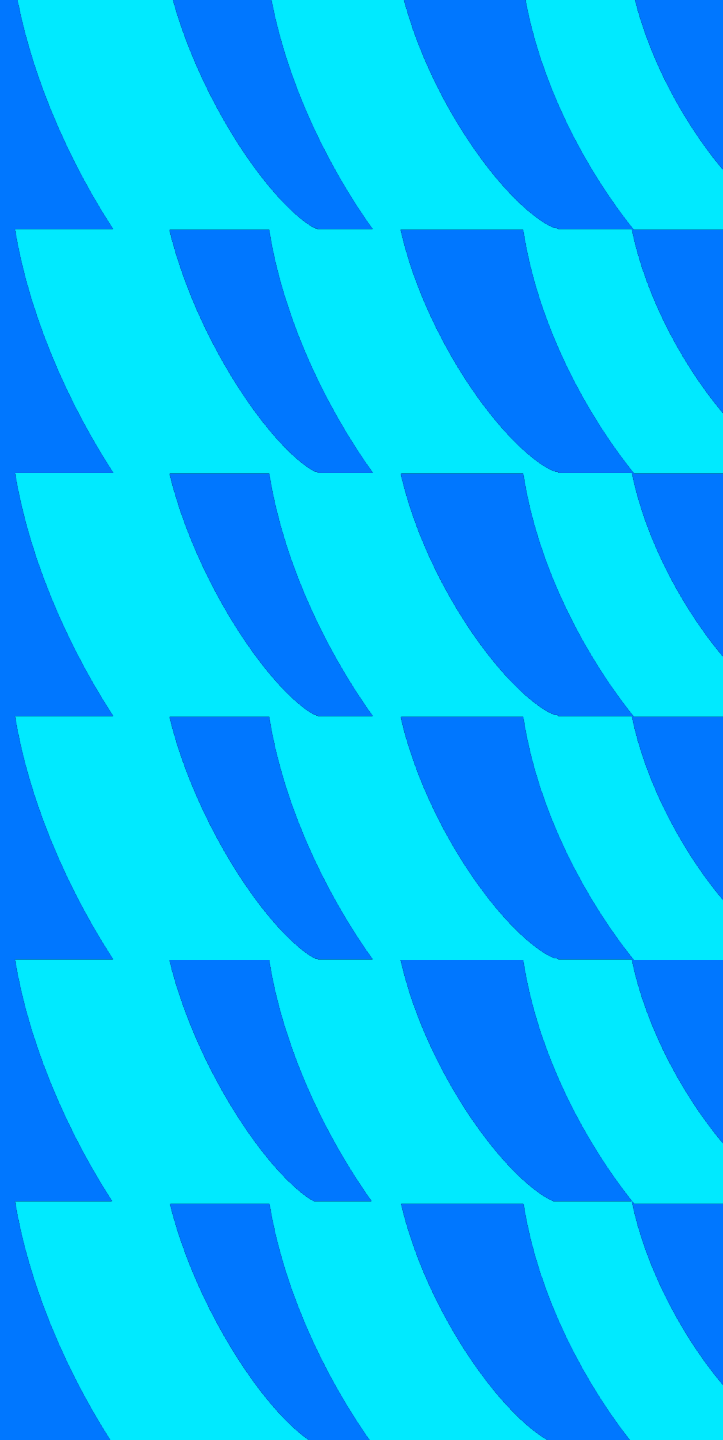


Обновили везде
критичные компоненты



Сделали устойчивый
к падениям кеширующий слой

Хозяйке
на заметку



Про **КОМПАКТНОСТЬ** объекта

Не забывай про выравнивание
и сортировку типов в объекте
Кеш скажет спасибо

AVX-инструкции часто могут дать выигрыш в несколько раз

Argon2 крайне чувствителен
к использованию векторных инструкций —
использование AVX-512F ускоряет
его работу в три раза по отношению
к обычной реализации



Переходы по указателям

НЕ БЕСПЛАТНЫ

1

Итерация в b-дереве
в 10x быстрее чем в set

2

Локальность обращений
к памяти очень важна

3

Заменяли стандартный
std::set на b-дерево

b-дерево вместо `std::set`

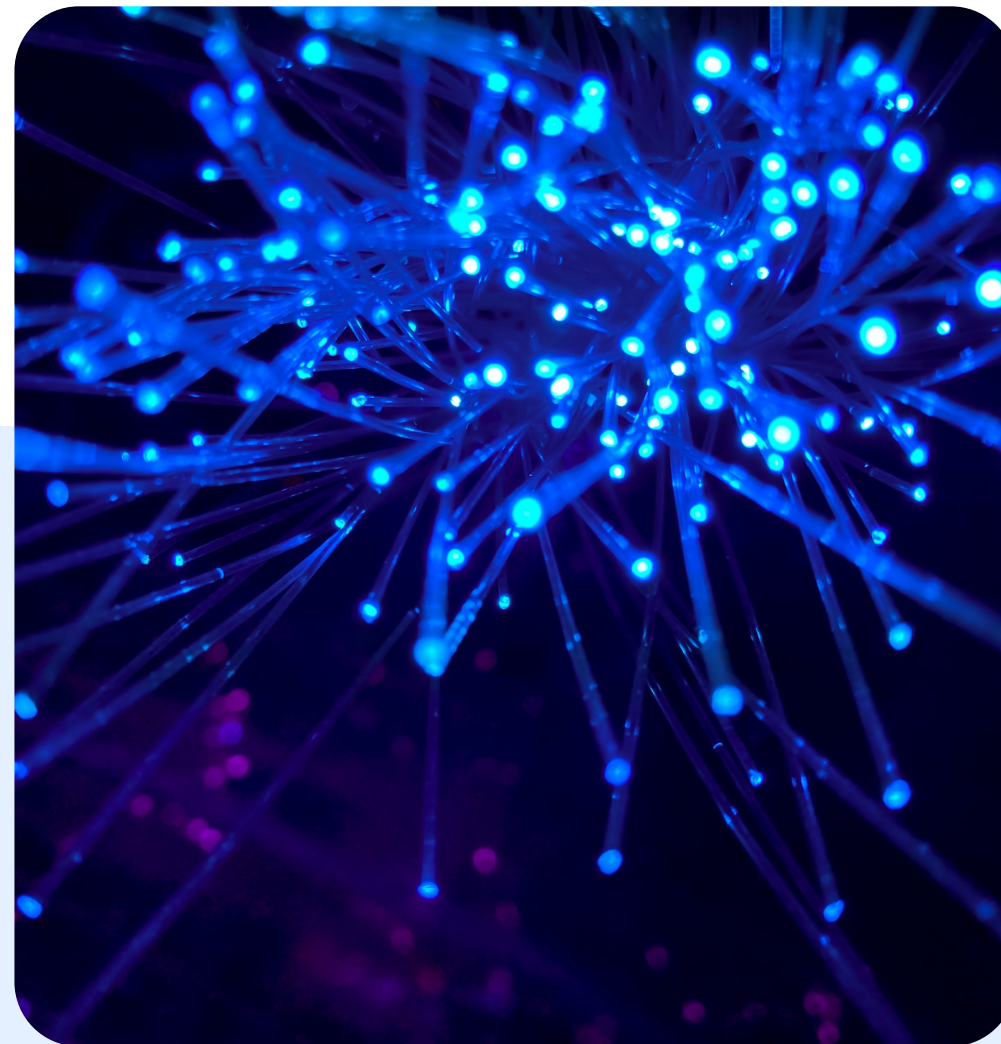
Вся служебная информация
в вершине b-дерева занимает 8 байт

- Размер поддереву вершины MAX 2^{31} [0:0, 3:0]
- Капасити вершины [3:1, 4,6]
- Фактический размер вершины [4:7, 6:2]
- Бит листа [6:3]



b-дерево вместо std::set || Профит

При факторе $B=32$ на один инстанс
мы тратим 4,8 байт, а сет, например,
40 байт



Используем устройство указателей для тегирования

Причем используем это двунаправленно,
младшие три бита указателя
из-за смещения также будут всегда нули

Старшие шестнадцать
для решения АВА-проблемы



Использование UDP-based- протокола в DC позволяет снизить задержку и сэкономить на CPU/RAM

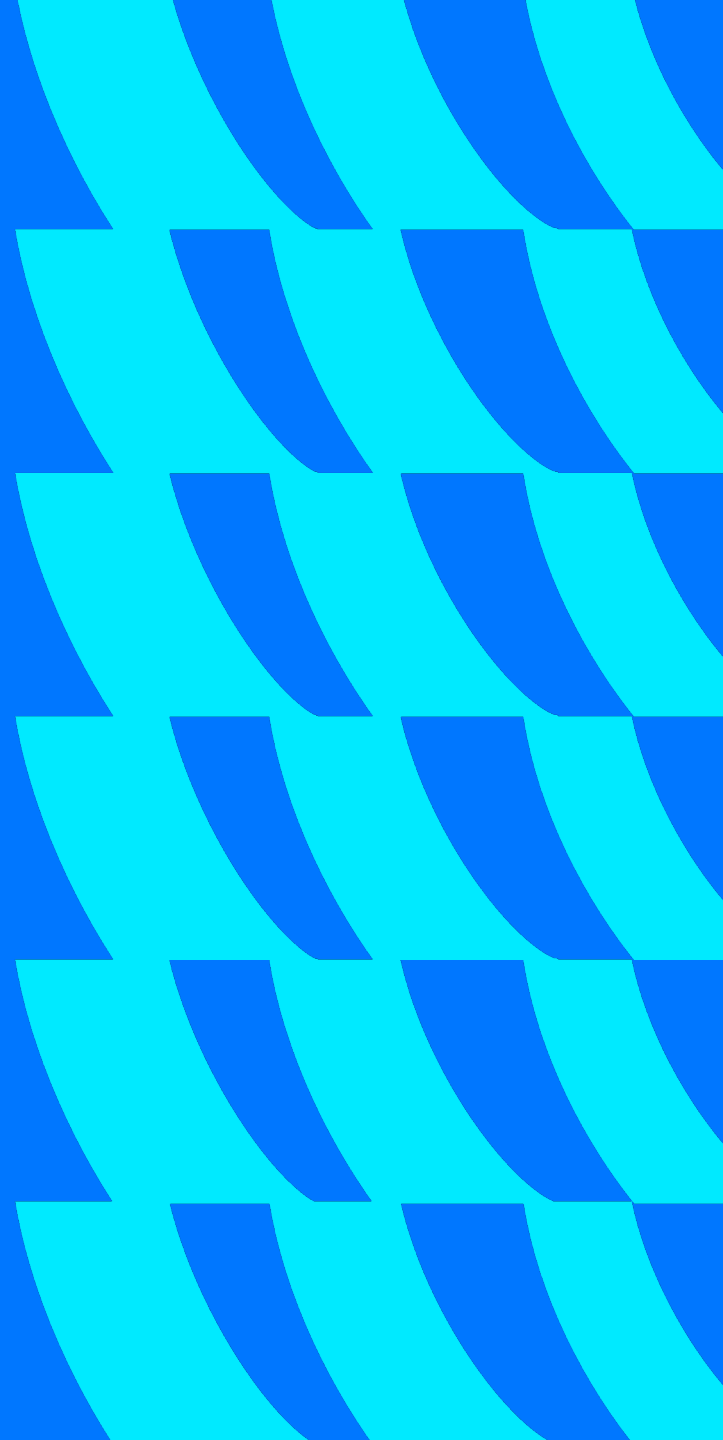
За счёт кастомной реализации
RPC-протокола на базе UDP
мы ушли от проблем:

- head of line blocking
- растущего хвоста latency
- лишне потребляемых ресурсов

Когда записи немного, стоит задуматься об RCU

Мы научились не копировать
всю структуру данных и отказались
от счетчиков ссылок для сущностей
Но об этом в другом докладе ;)

Вместо заключения





Бенчмарк-
бенчмарк-бенчмарк,
не доверяйте
экспертам



Используйте фазеры
и property-based-
тестирование



Будьте аккуратны
с продакшном



Ошибаются все,
главное – правильно
сделать выводы



Спасибо
за внимание!